FLEXIBILITY

EFFICIENCY

CPUs

Image: Microsoft

Image: [Microsoft](#)

Image: Microsoft

CPUs     GPUs     FPGAs     ASICs

FLEXIBILITY                 EFFICIENCY

Image: Microsoft

First stage of
LHC *trigger*

CPUs

GPUs

FPGAs

ASICs

FLEXIBILITY

EFFICIENCY

Image: Microsoft
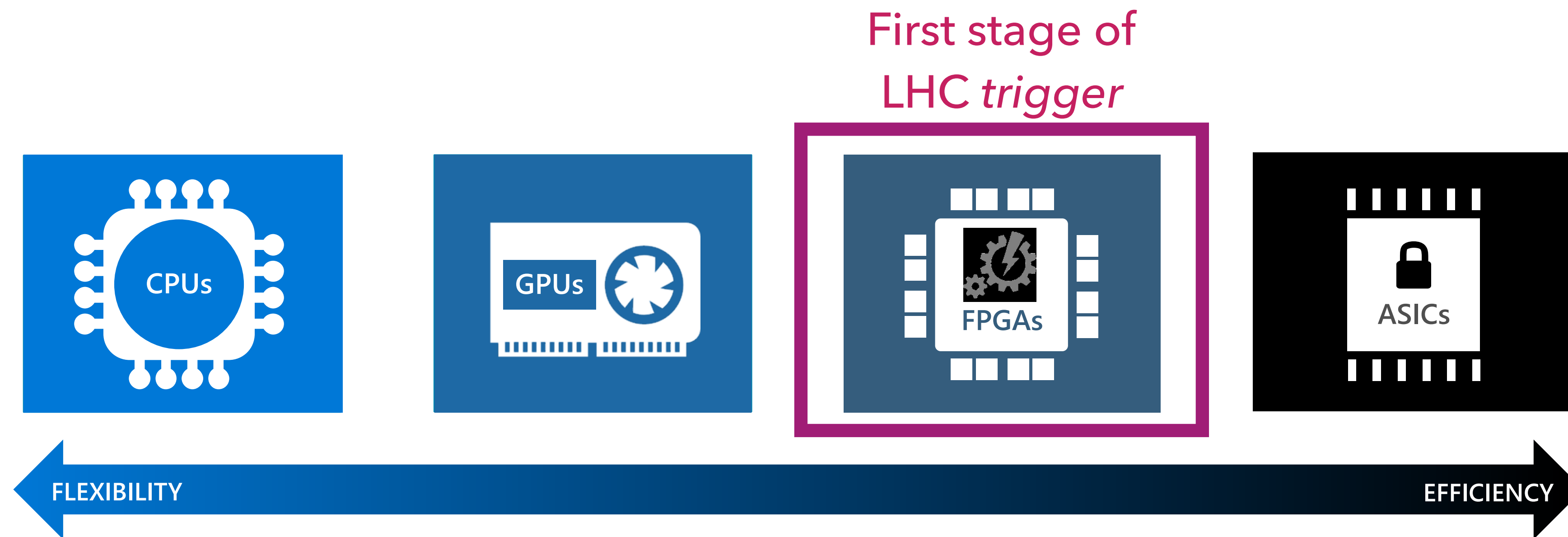
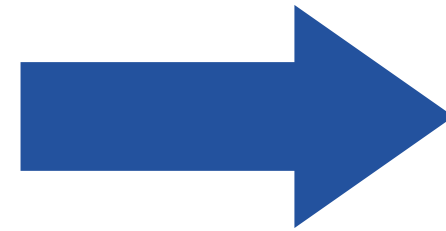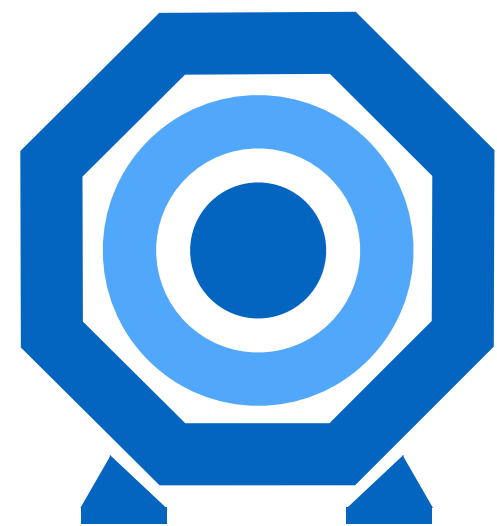Compute
Latency

1 ns                    1 μs                    1 ms                    1 s

**40 MHz**

ASICs

**Challenges:**

Each collision produces $O(10^3)$ particles

The detectors have $O(10^8)$ sensors
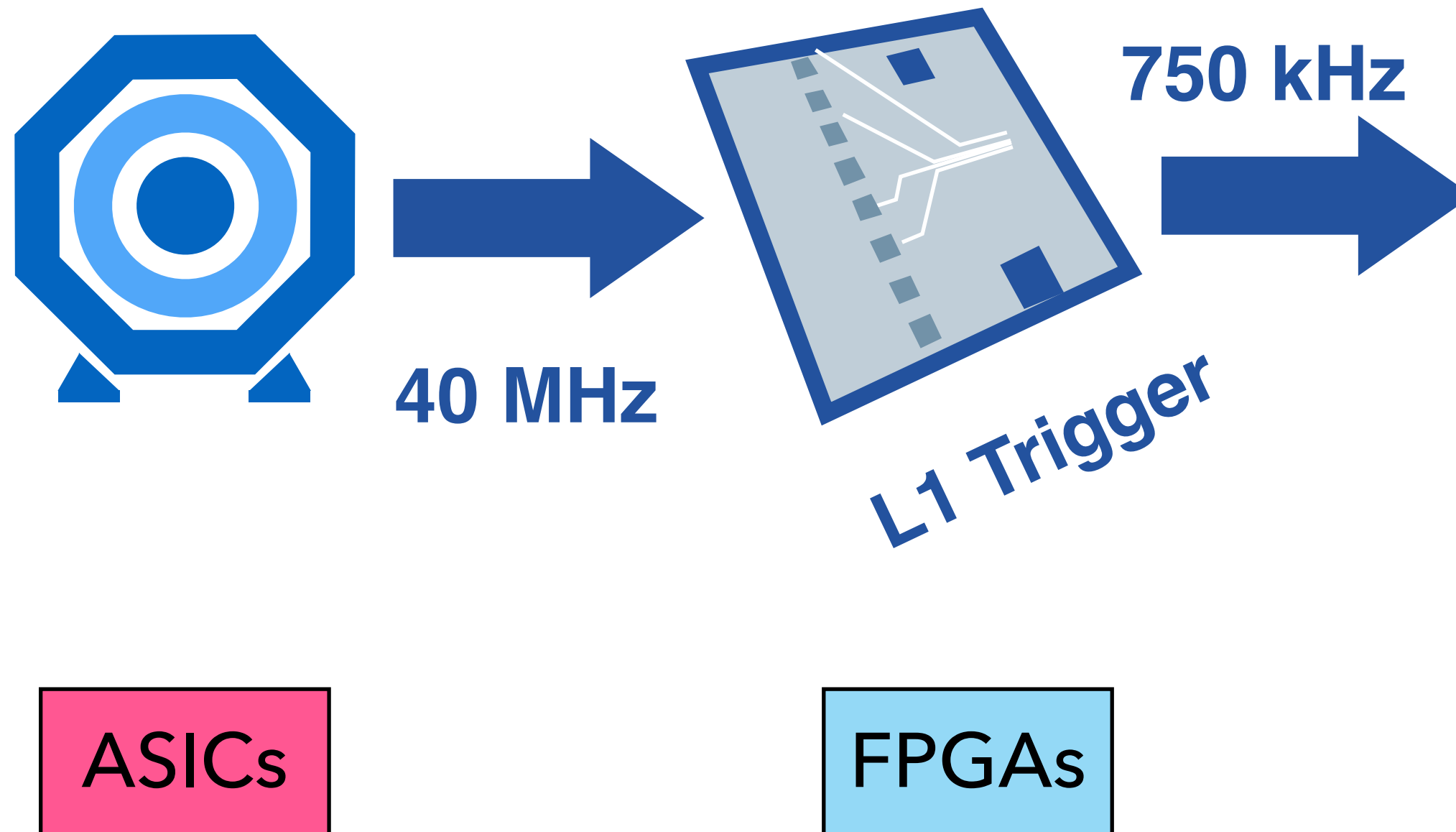
Extreme data rates of $O(100 \text{ TB/s})$

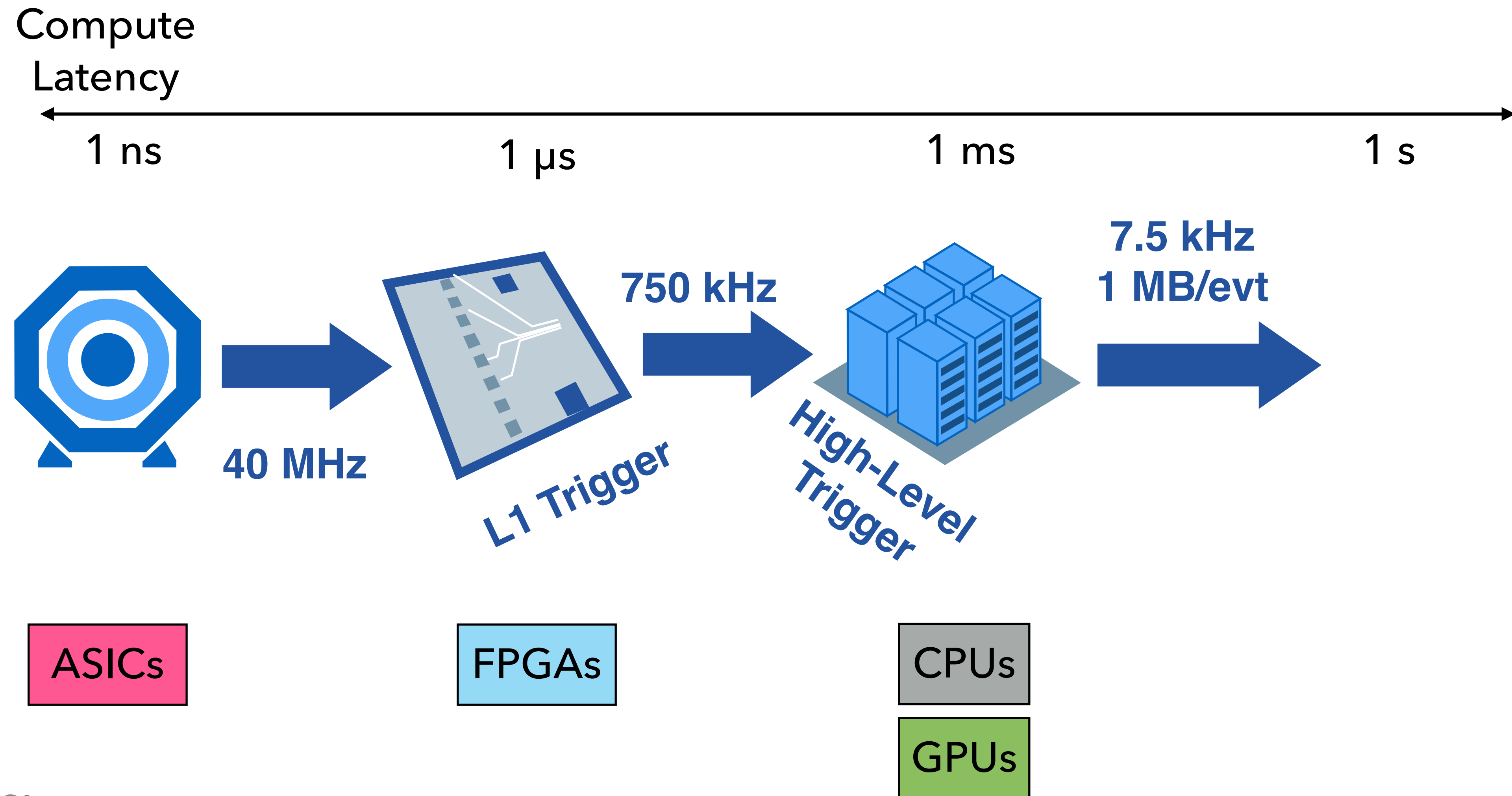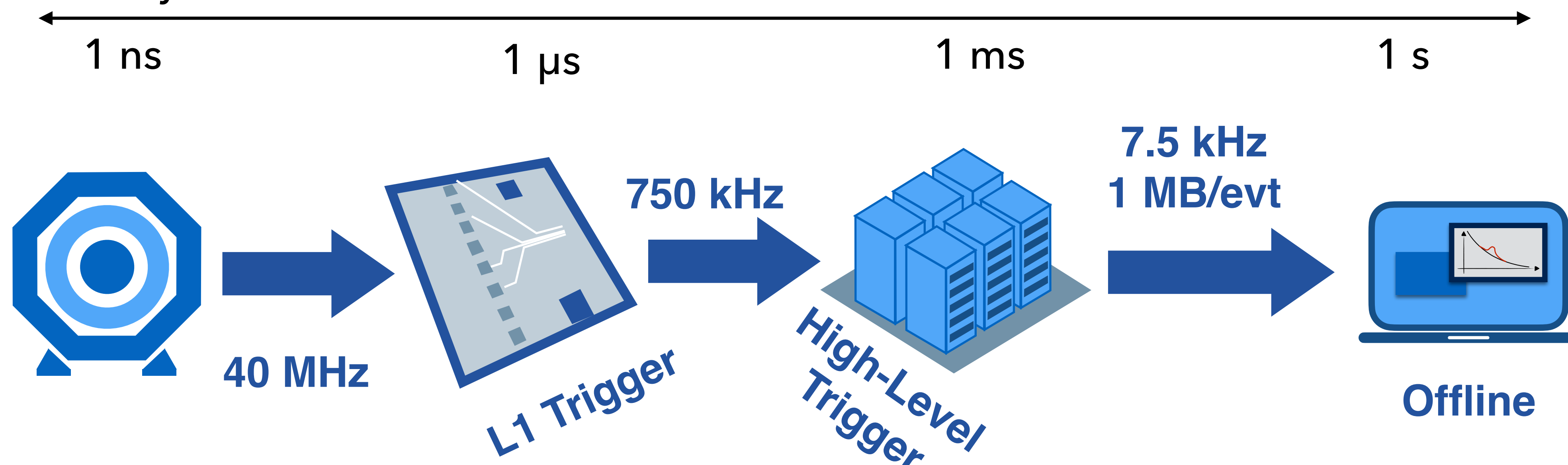Compute
Latency

1 ns          1 μs          1 ms          1 s

40 MHz

750 kHz

L1 Trigger

ASICs          FPGAs

**Challenges:**

Each collision produces O($10^3$) particles

The detectors have O($10^8$) sensors

Extreme data rates of O(100 TB/s)

Compute
Latency

1 ns          1 µs          1 ms          1 s



40 MHz    L1 Trigger    750 kHz    High-Level Trigger    7.5 kHz
1 MB/evt

ASICs          FPGAs          CPUs

GPUs

**Challenges:**

Each collision produces O($10^3$) particles

The detectors have O($10^8$) sensors

Extreme data rates of O(100 TB/s)

Compute Latency

1 ns       1 μs       1 ms       1 s

40 MHz    **L1 Trigger**    **750 kHz**    **High-Level Trigger**    **7.5 kHz 1 MB/evt**    **Offline**

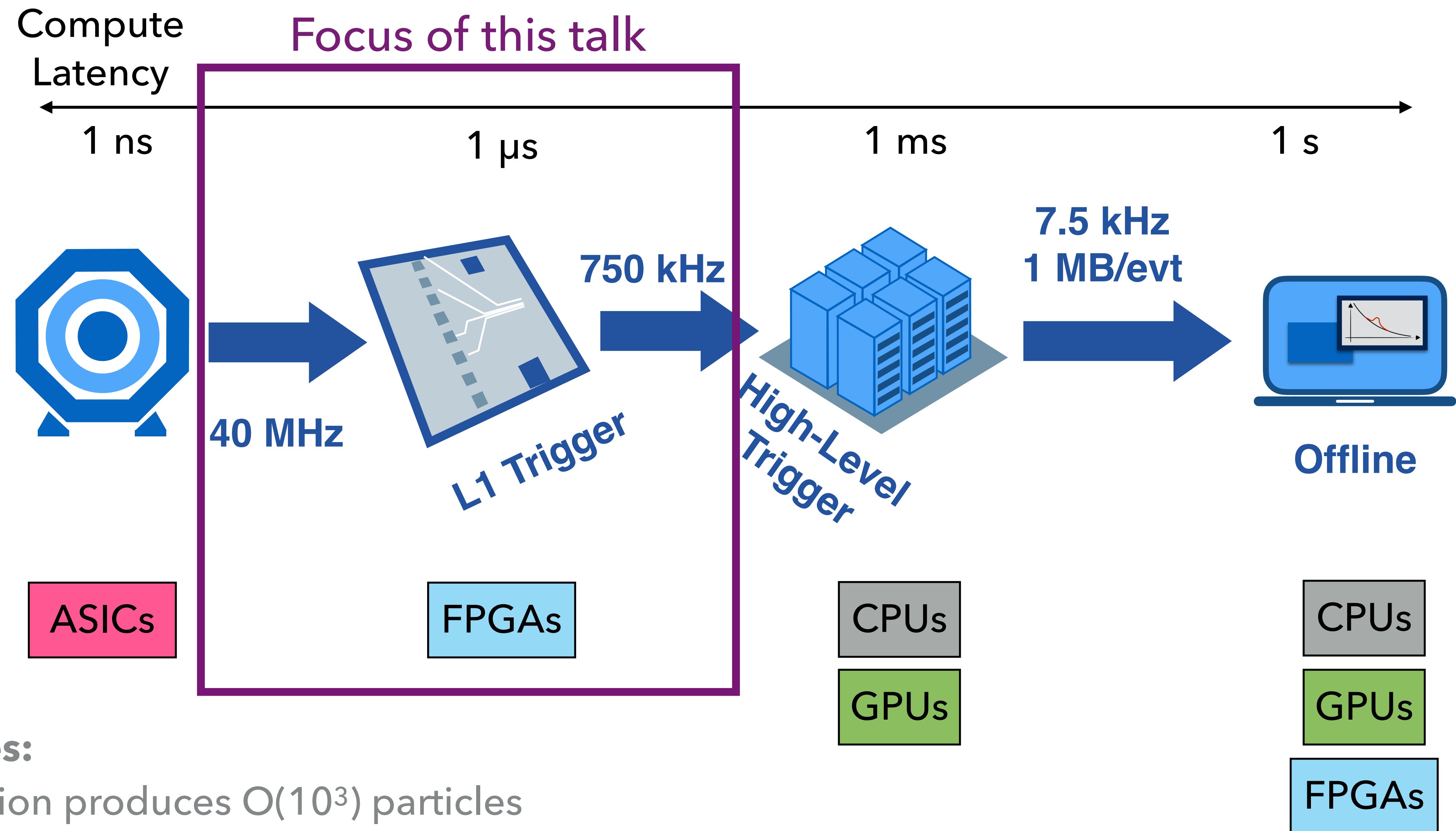ASICs     FPGAs     CPUs     CPUs

GPUs     GPUs

FPGAs

**Challenges:**

Each collision produces $O(10^3)$ particles

The detectors have $O(10^8)$ sensors

Extreme data rates of $O(100\ \text{TB/s})$

Exabyte-scale datasets

Compute Latency

Focus of this talk

1 ns       1 μs       1 ms       1 s

40 MHz

L1 Trigger

750 kHz

High-Level Trigger

7.5 kHz
1 MB/evt

Offline

ASICs

FPGAs

CPUs

GPUs

CPUs

GPUs

FPGAs

**Challenges:**

Each collision produces $O(10^3)$ particles

The detectors have $O(10^8)$ sensors

Extreme data rates of $O(100\ TB/s)$

Exabyte-scale datasets

Thresholds set by
backgrounds,
limited resolution
@ L1, and rate
budget

CMS-TDR-021

| Trigger | Threshold [GeV] | 6 |
| --- | --- | --- |

Thresholds set by
backgrounds,
limited resolution
@ L1, and rate
budget

CMS-TDR-021

# SIMPLIFIED HL–LHC L1 TRIGGER MENU

▸ Single/double/triple muons/electrons

| Trigger | Threshold [GeV] |
|---------|-----------------|
| 1 μ | 22 |
| 2 μ | 15, 7 |
| 3 μ | 5, 3, 3 |
| 1 e | 36 |
| 2 e | 25, 12 |

Thresholds set by backgrounds, limited resolution @ L1, and rate budget

CMS-TDR-021

# SIMPLIFIED HL–LHC L1 TRIGGER MENU

▸ Single/double/triple muons/electrons

▸ Photons

| Trigger | Threshold [GeV] |
|---------|-----------------|
| 1 μ | 22 |
| 2 μ | 15, 7 |
| 3 μ | 5, 3, 3 |
| 1 e | 36 |
| 2 e | 25, 12 |
| 1 γ | 36 |
| 2 γ | 22, 12 |

Thresholds set by backgrounds, limited resolution @ L1, and rate budget
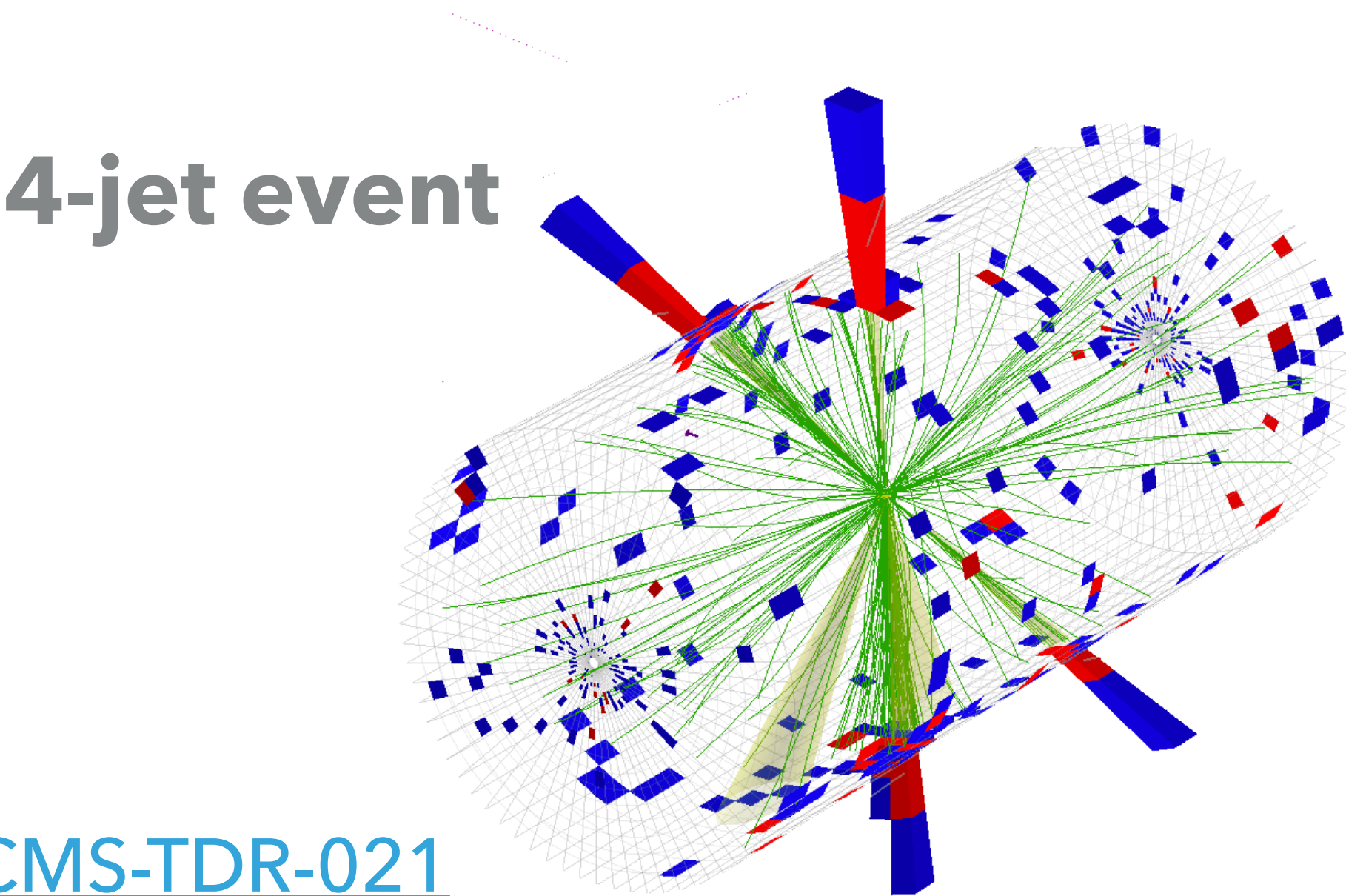
# SIMPLIFIED HL–LHC L1 TRIGGER MENU

▸ Single/double/triple muons/electrons

▸ Photons

▸ Taus

| Trigger | Threshold [GeV] |
|---------|-----------------|
| 1 μ | 22 |
| 2 μ | 15, 7 |
| 3 μ | 5, 3, 3 |
| 1 e | 36 |
| 2 e | 25, 12 |
| 1 γ | 36 |
| 2 γ | 22, 12 |
| 1 τ | 150 |
| 2 τ | 90, 90 |

Thresholds set by backgrounds, limited resolution @ L1, and rate budget

# SIMPLIFIED HL–LHC L1 TRIGGER MENU

- ▸ Single/double/triple muons/electrons
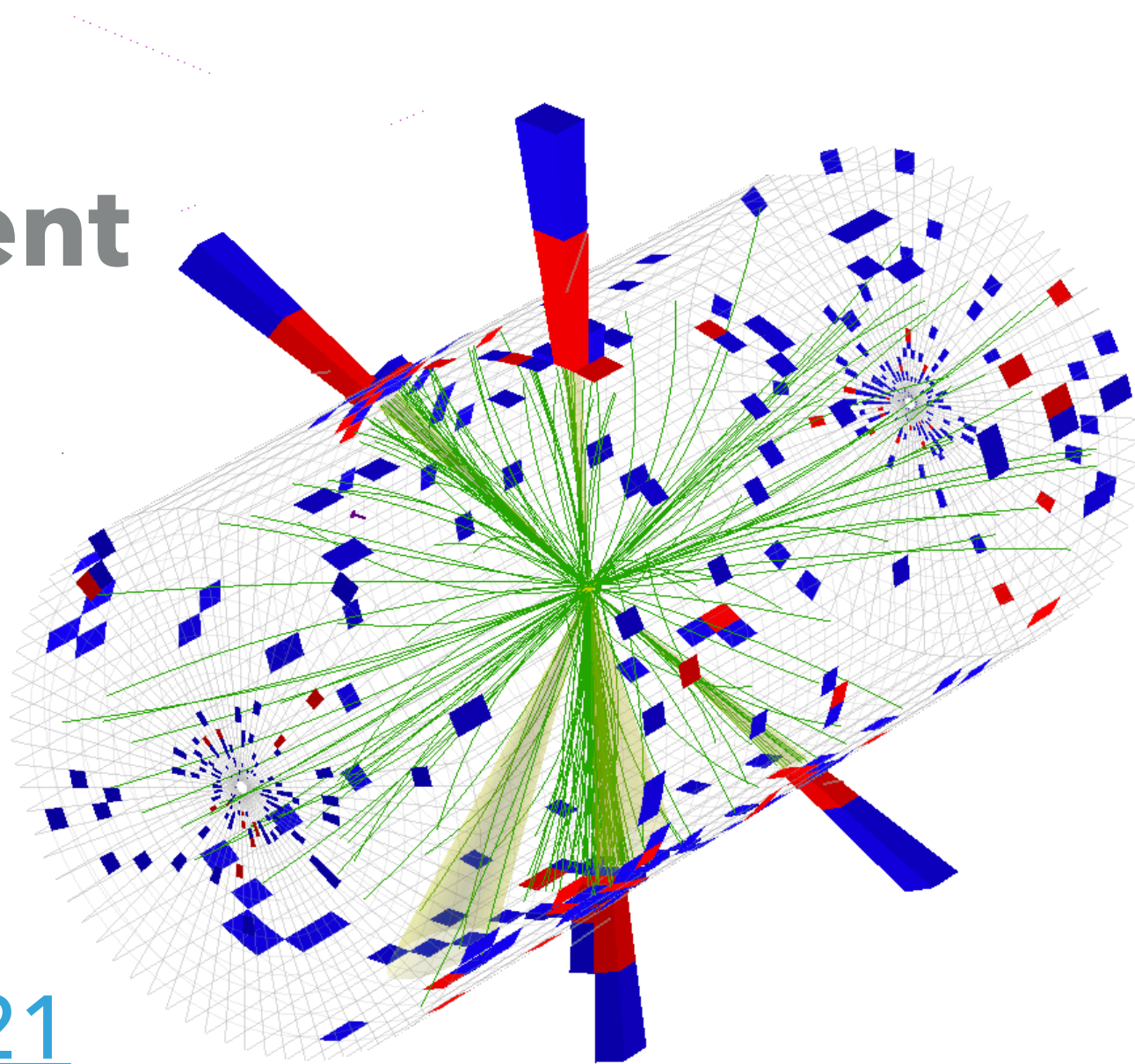- ▸ Photons
- ▸ Taus
- ▸ Hadronic

**4-jet event**



Thresholds set by backgrounds, limited resolution @ L1, and rate budget

CMS-TDR-021

| Trigger | Threshold [GeV] |
|---|---|
| 1 μ | 22 |
| 2 μ | 15, 7 |
| 3 μ | 5, 3, 3 |
| 1 e | 36 |
| 2 e | 25, 12 |
| 1 γ | 36 |
| 2 γ | 22, 12 |
| 1 τ | 150 |
| 2 τ | 90, 90 |
| 1 jet | 180 |
| 2 jet | 112, 112 |
| $H_T$ | 450 |
| 4 jet + $H_T$ | 75, 55, 40, 40, 400 |

# SIMPLIFIED HL–LHC L1 TRIGGER MENU

▸ Single/double/triple muons/electrons

▸ Photons

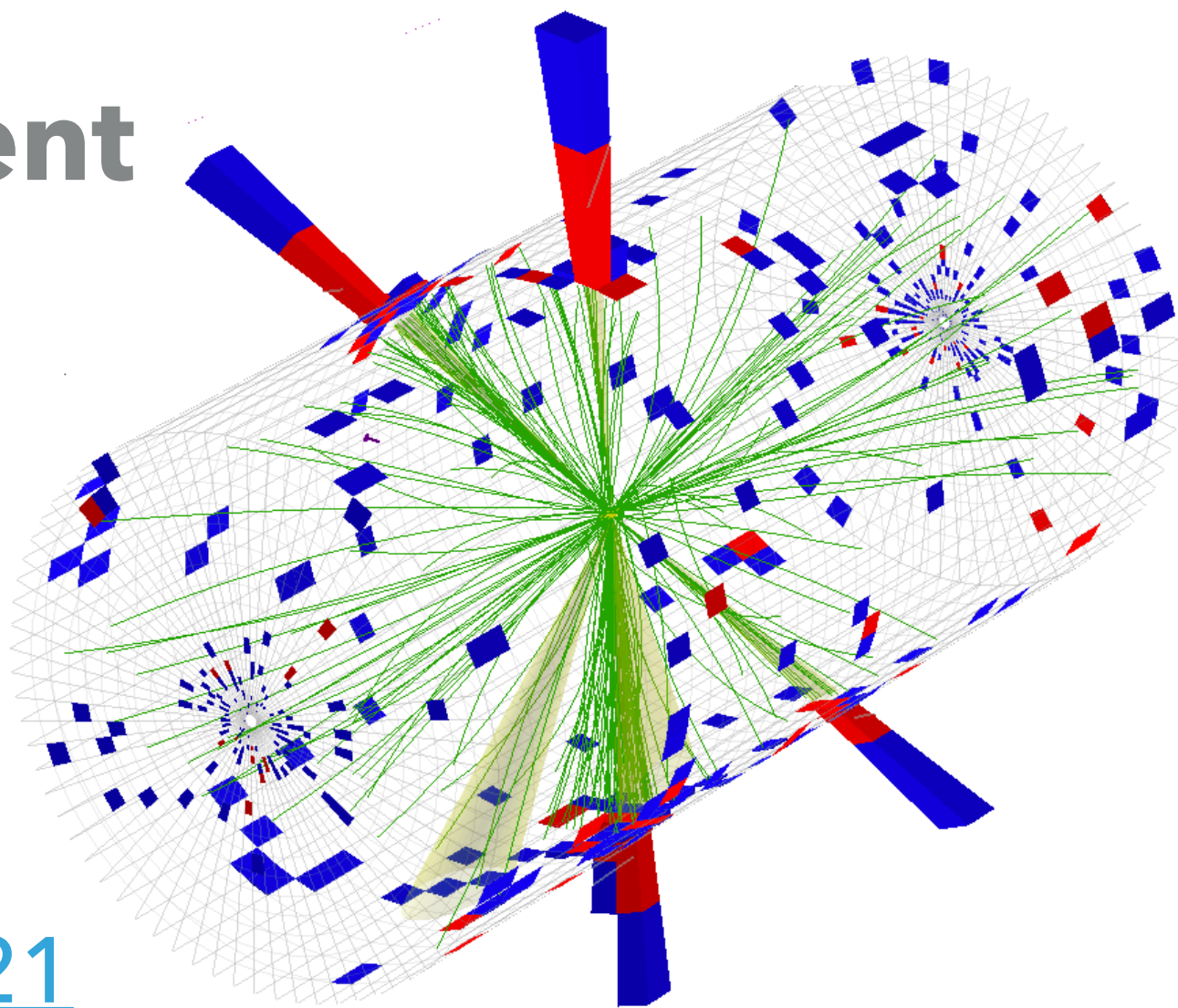▸ Taus

▸ Hadronic

▸ Missing transverse energy

**4-jet event**



Thresholds set by backgrounds, limited resolution @ L1, and rate budget

CMS-TDR-021

| Trigger | Threshold [GeV] |
|---------|-----------------|
| 1 μ | 22 |
| 2 μ | 15, 7 |
| 3 μ | 5, 3, 3 |
| 1 e | 36 |
| 2 e | 25, 12 |
| 1 γ | 36 |
| 2 γ | 22, 12 |
| 1 τ | 150 |
| 2 τ | 90, 90 |
| 1 jet | 180 |
| 2 jet | 112, 112 |
| $H_T$ | 450 |
| 4 jet + $H_T$ | 75, 55, 40, 40, 400 |
| $p_T^{miss}$ | 200 |

# SIMPLIFIED HL-LHC L1 TRIGGER MENU

▸ Single/double/triple muons/electrons

▸ Photons

▸ Taus

▸ Hadronic

▸ Missing transverse energy

▸ "Cross" triggers (not shown)

**4-jet event**

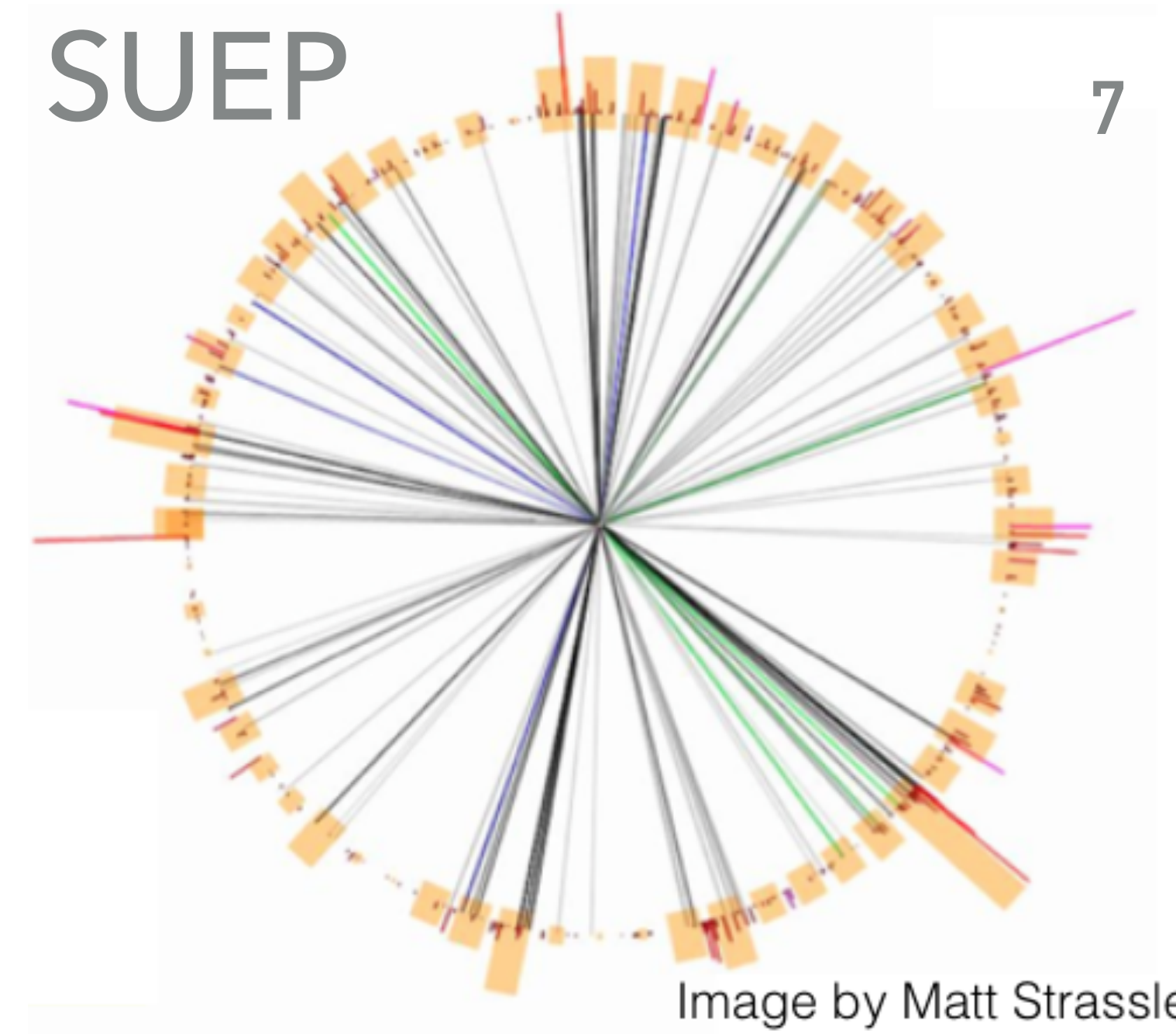Thresholds set by backgrounds, limited resolution @ L1, and rate budget

CMS-TDR-021

| Trigger | Threshold [GeV] |
|---|---|
| 1 μ | 22 |
| 2 μ | 15, 7 |
| 3 μ | 5, 3, 3 |
| 1 e | 36 |
| 2 e | 25, 12 |
| 1 γ | 36 |
| 2 γ | 22, 12 |
| 1 τ | 150 |
| 2 τ | 90, 90 |
| 1 jet | 180 |
| 2 jet | 112, 112 |
| $H_T$ | 450 |
| 4 jet + $H_T$ | 75, 55, 40, 40, 400 |
| $p_T^{miss}$ | 200 |

# WHAT COULD WE BE MISSING?

▸ How can we trigger on more complex low-energy hadronic signatures? Long-lived/displaced particles?
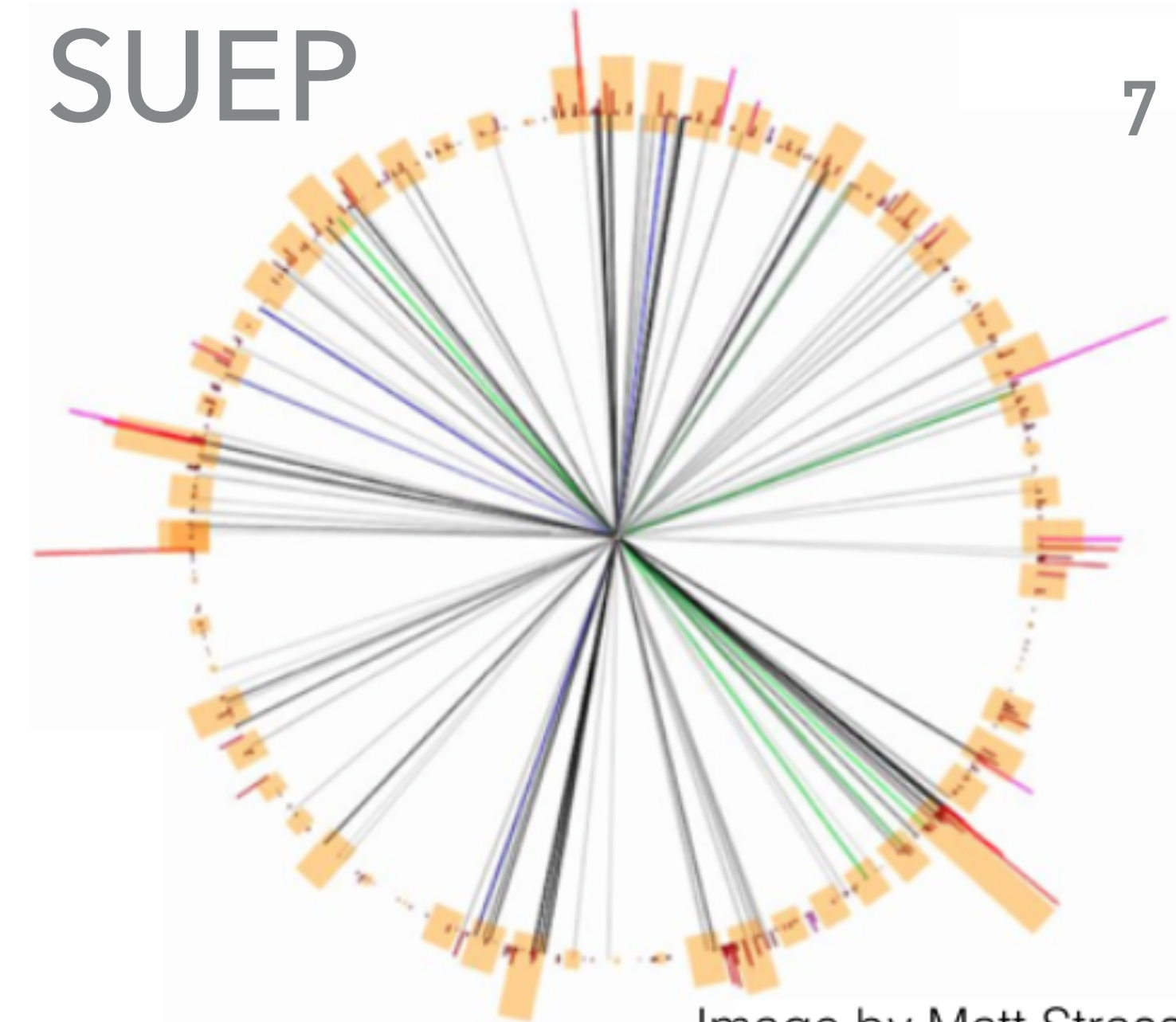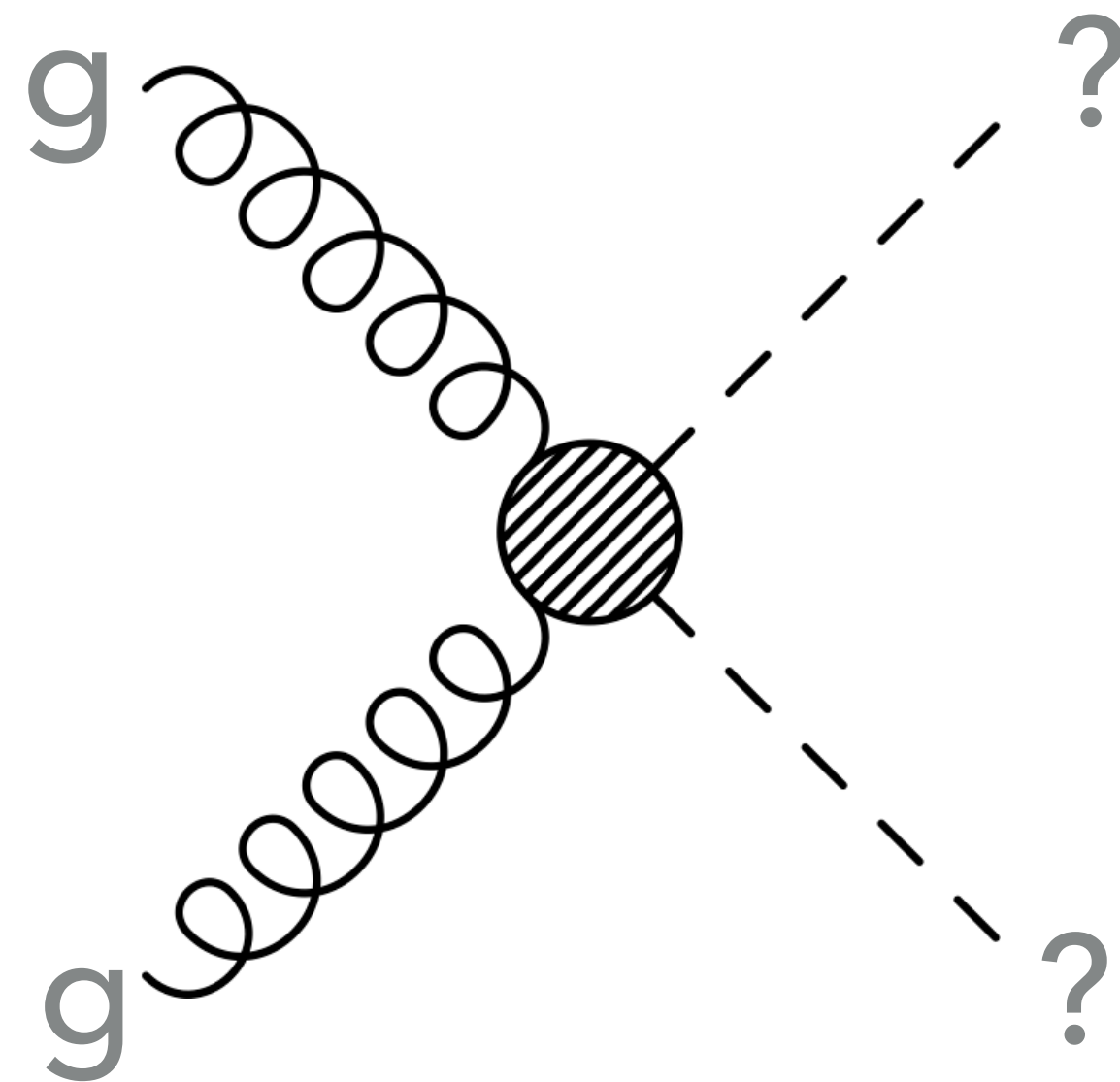
SUEP



Image by Matt Strassler

# WHAT COULD WE BE MISSING?

- ▸ How can we trigger on more complex low-energy hadronic signatures? Long-lived/displaced particles?
- ▸ What if we don't know exactly what to look for?

Image by Matt Strassler

# WHAT COULD WE BE MISSING?

▸ How can we trigger on more complex low-energy hadronic signatures? Long-lived/displaced particles?

▸ What if we don't know exactly what to look for?

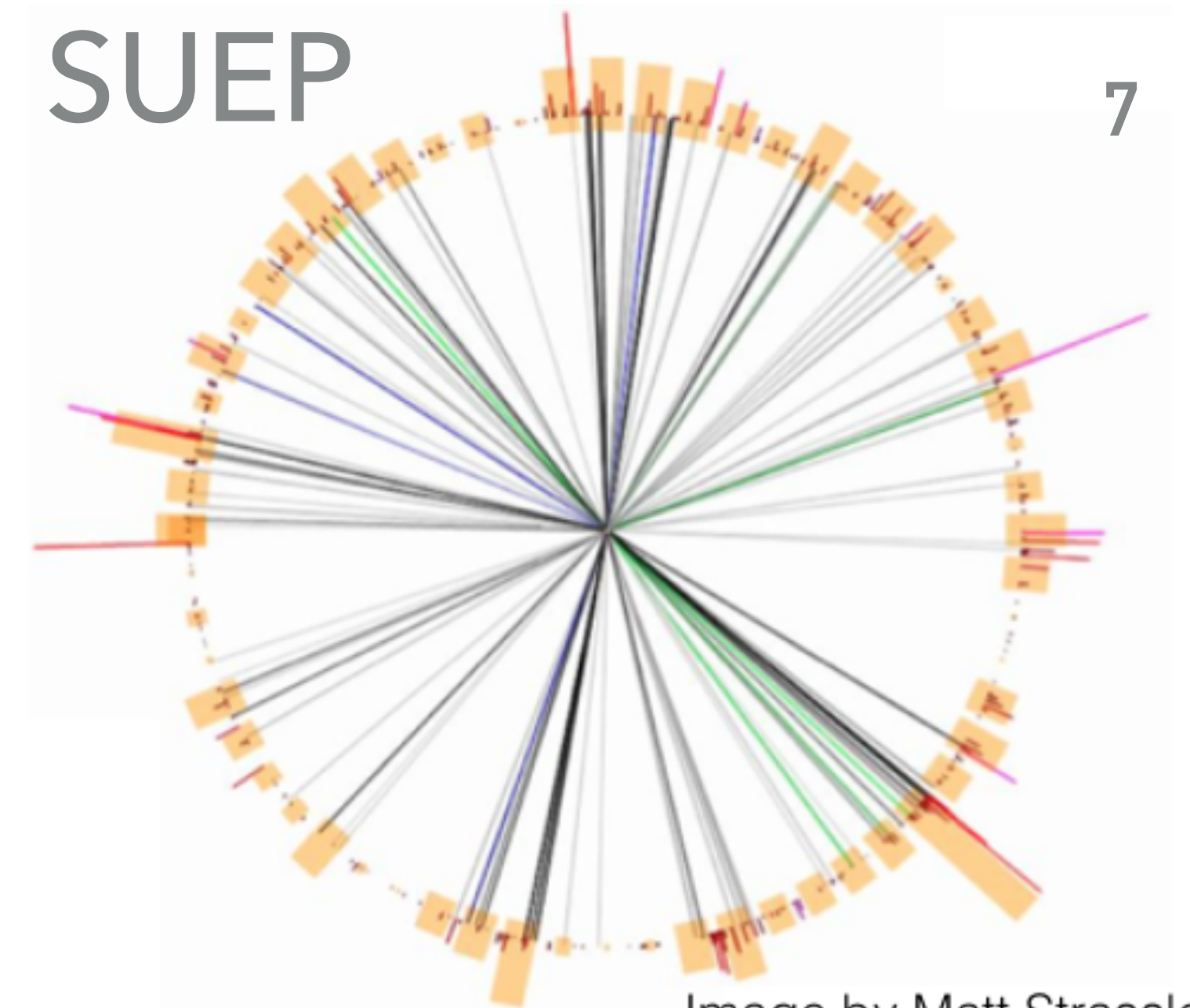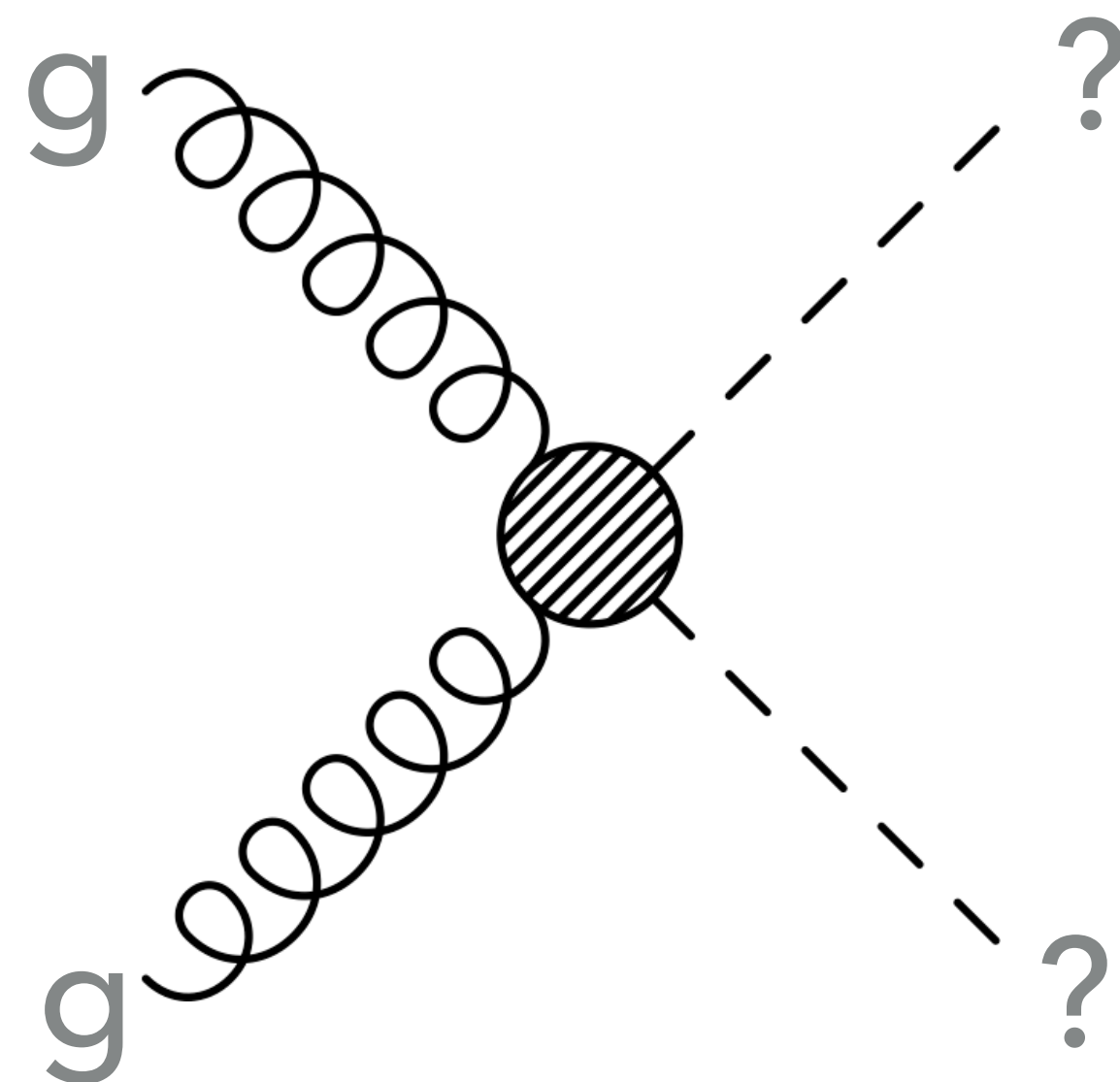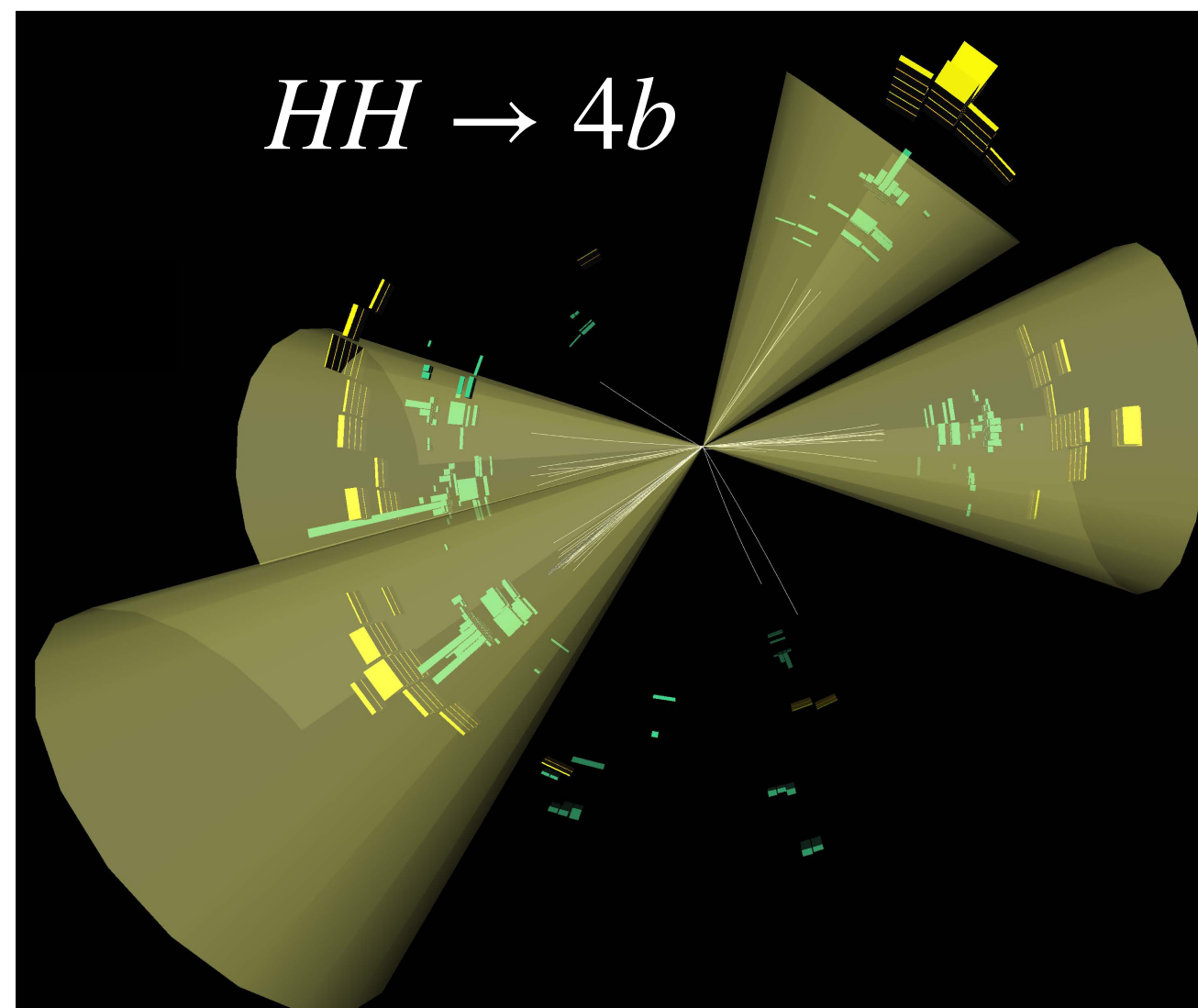▸ What if our signatures require complex multivariate algorithms (e.g. b tagging)?
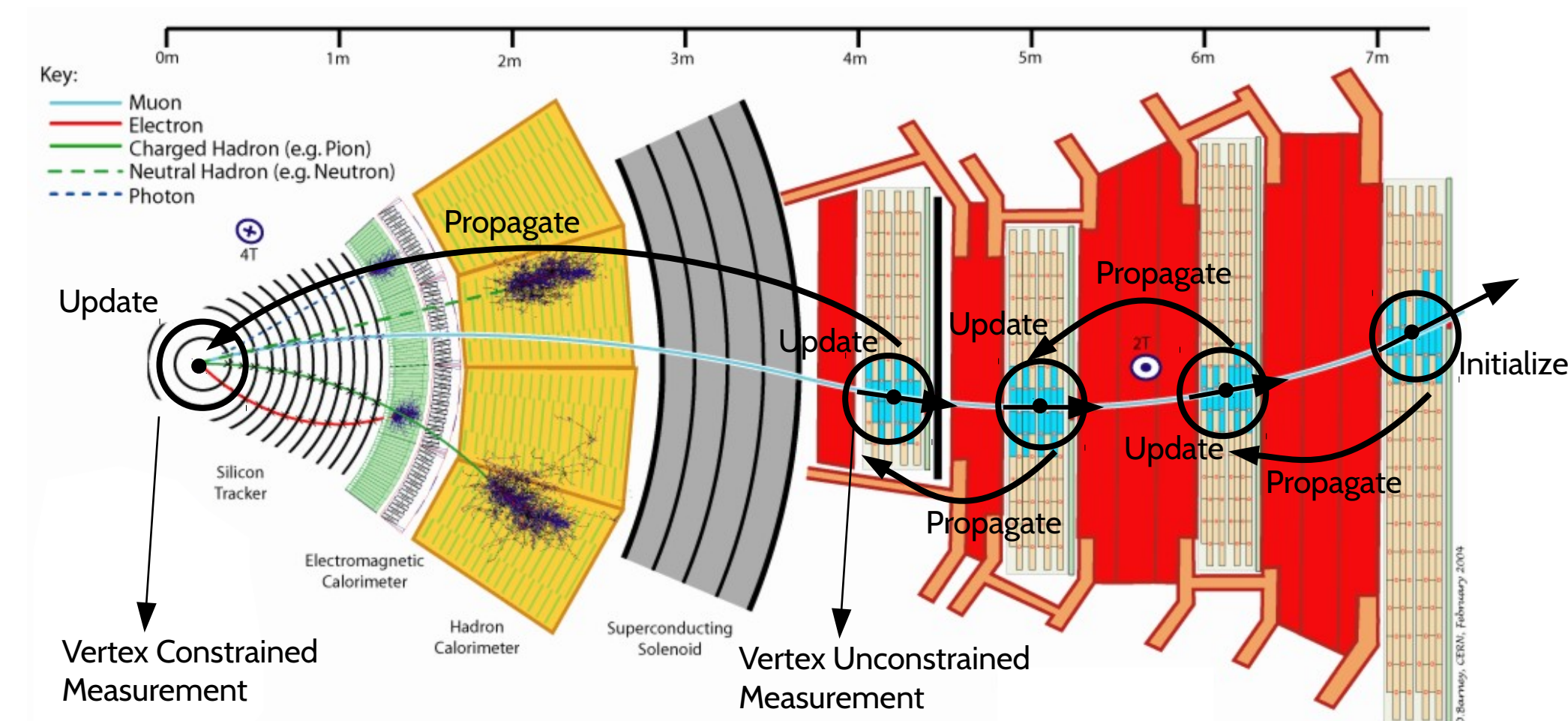


Image by Matt Strassler



$HH \rightarrow 4b$

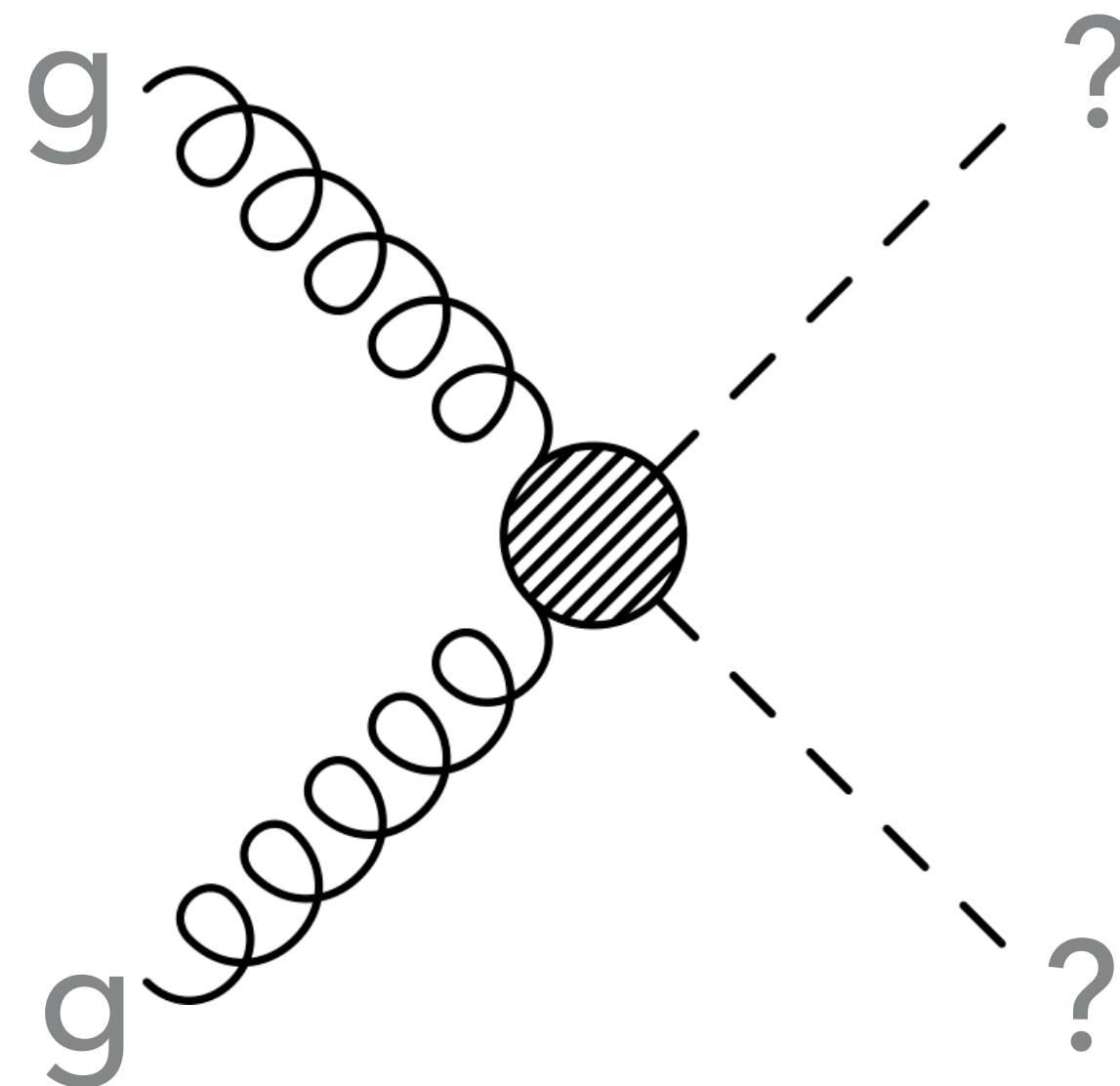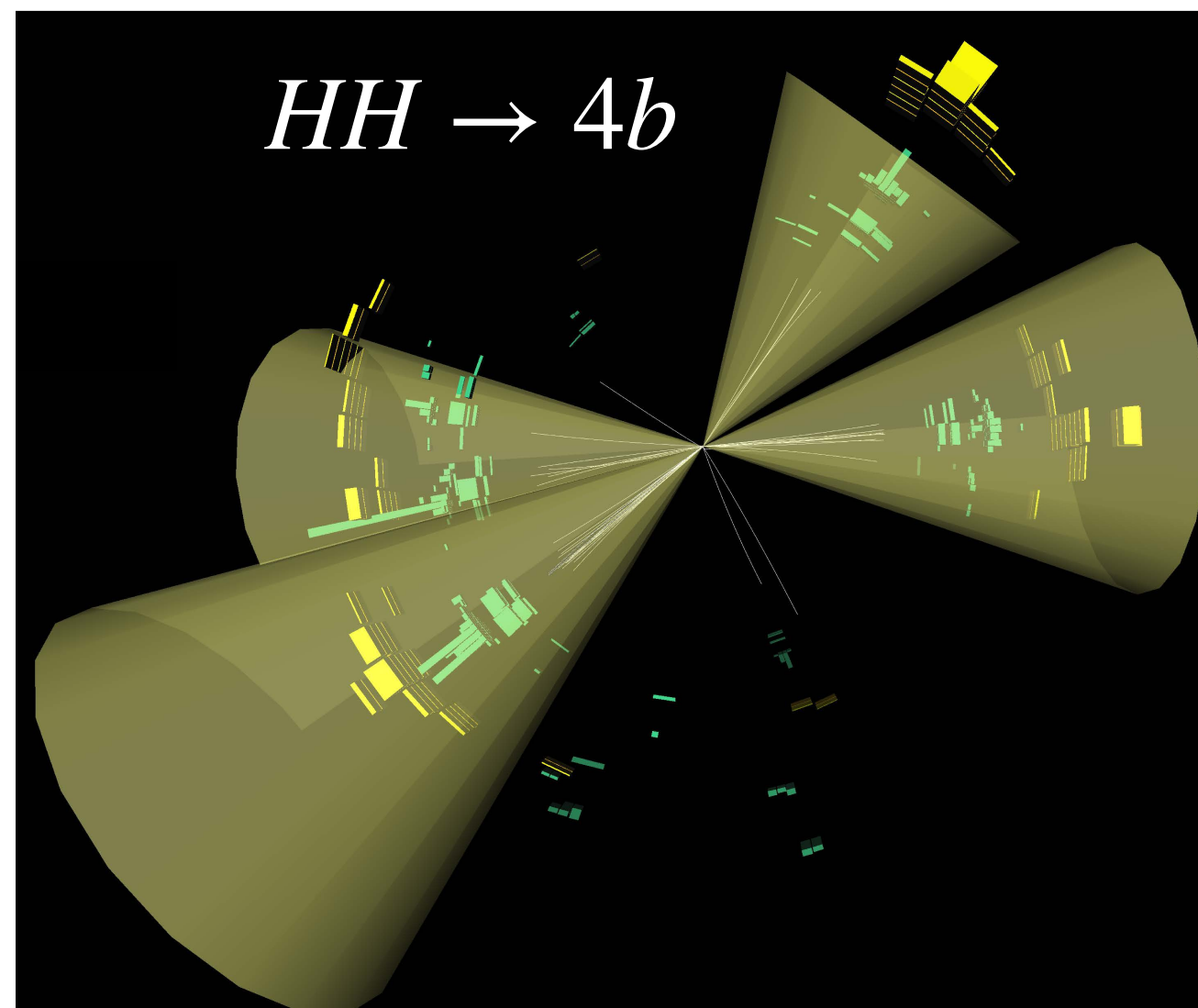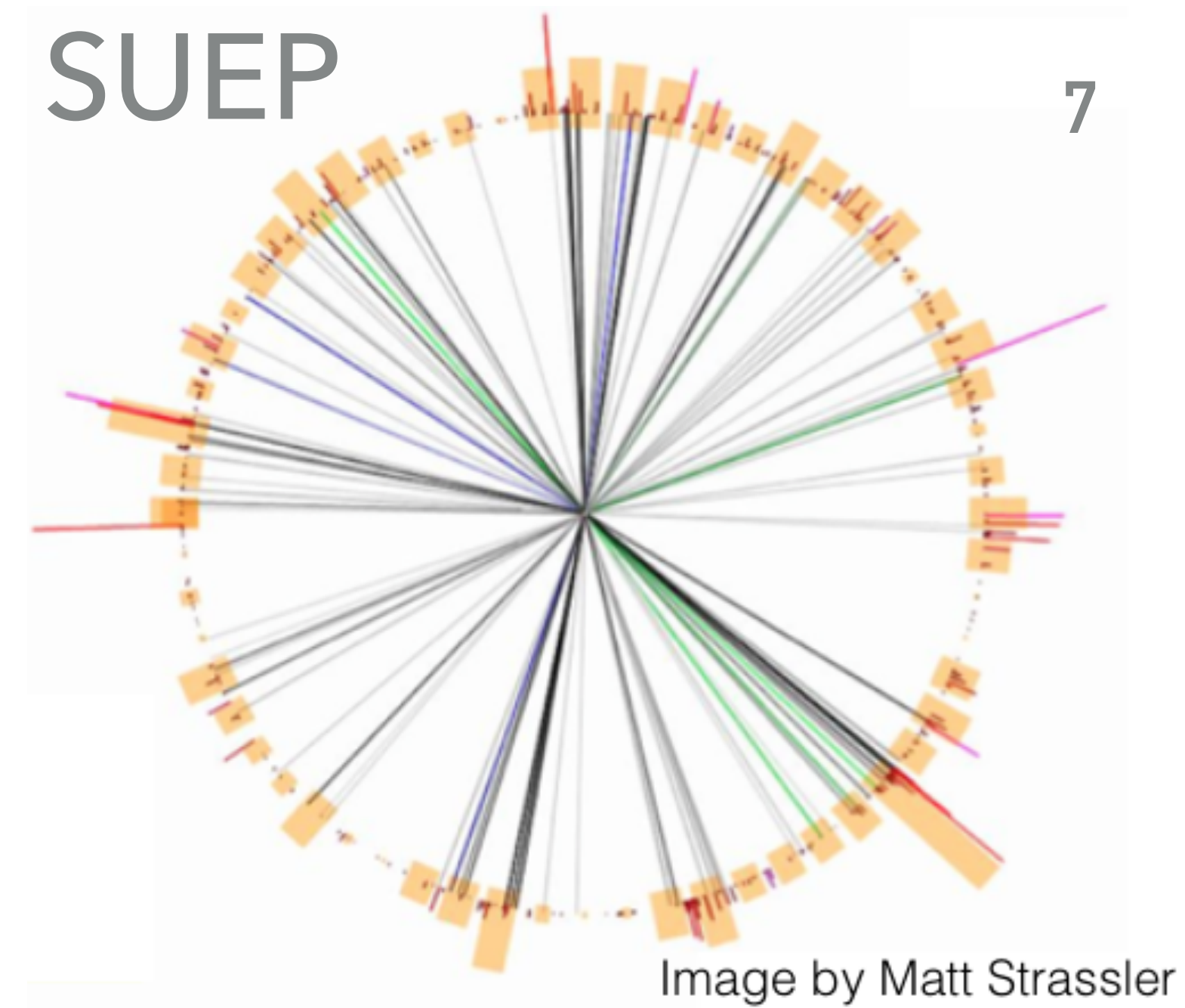# WHAT COULD WE BE MISSING?

▸ How can we trigger on more complex low-energy hadronic signatures? Long-lived/displaced particles?

▸ What if we don't know exactly what to look for?

▸ What if our signatures require complex multivariate algorithms (e.g. b tagging)?

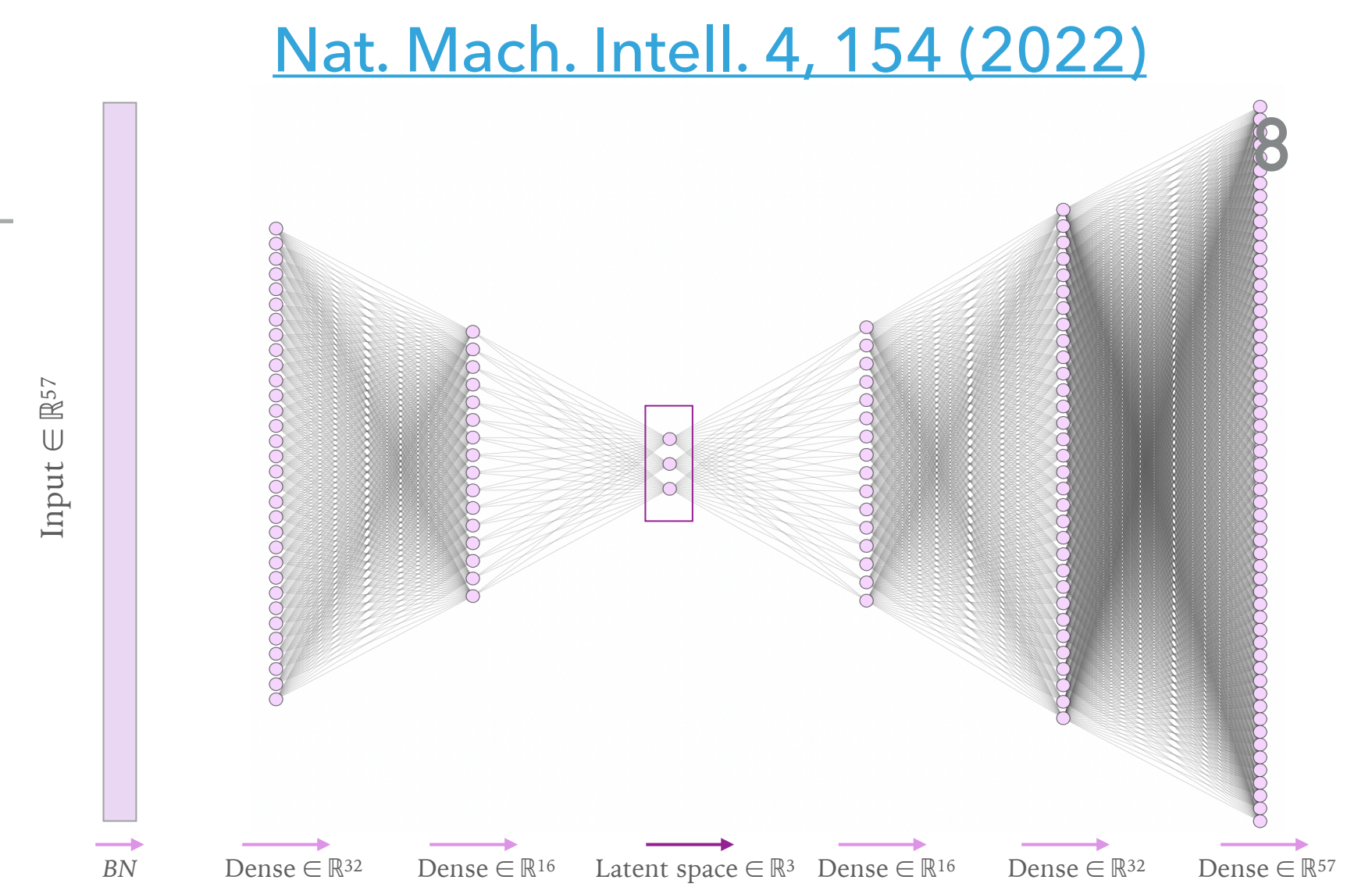▸ How can we improve on our traditional (often slow) reconstruction algorithms?
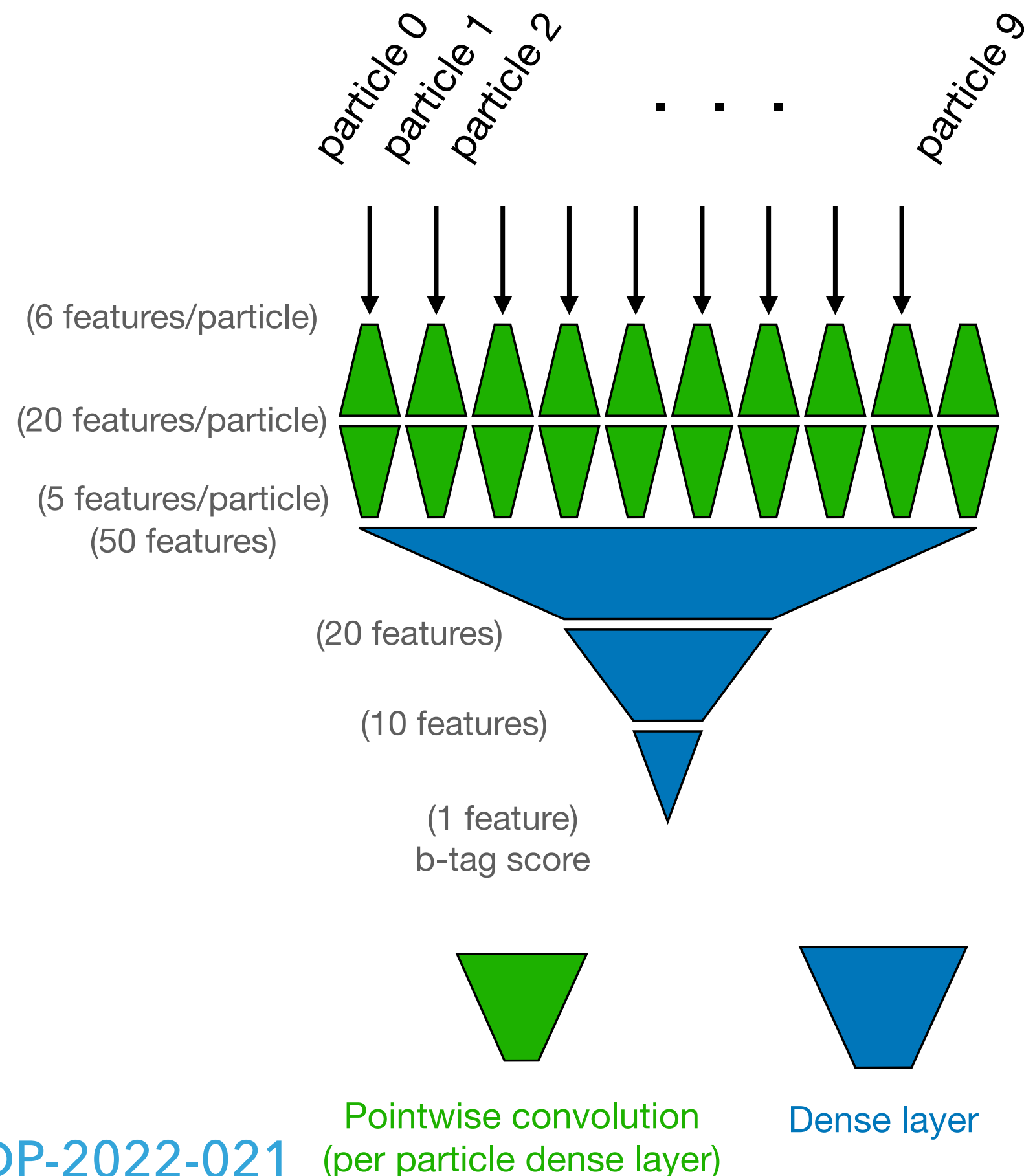
## SUEP



Image by Matt Strassler



$HH \rightarrow 4b$



g

g

?

?

# ML IN THE TRIGGER

▶ (Variational) autoencoders for anomaly detection

Input $\in \mathbb{R}^{57}$

$BN$    Dense $\in \mathbb{R}^{32}$    Dense $\in \mathbb{R}^{16}$    Latent space $\in \mathbb{R}^3$    Dense $\in \mathbb{R}^{16}$    Dense $\in \mathbb{R}^{32}$    Dense $\in \mathbb{R}^{57}$

Block 1:
Conv2d (16,(3,3))
ReLU
AvPooling (3,1)

Block 2:
Conv2d 1 (32,(3,1))
ReLU
AvPooling (3,1)
Flatten (64)

Block 3:
Dense (8)
Dense 1 (64)
ReLU
Reshape (2,1,32)

Block 4:
Conv2d 2 (32,(3,1))
ReLU
UpSampling (3,1)
ZeroPad (0,0),(1,1)

Block 5:
Conv2d 3 (16,(3,1))
ReLU
UpSampling (3,1)
ZeroPad (1,0),(0,0)

Block 0:
Input 19x3x1
ZeroPadding (1,0)
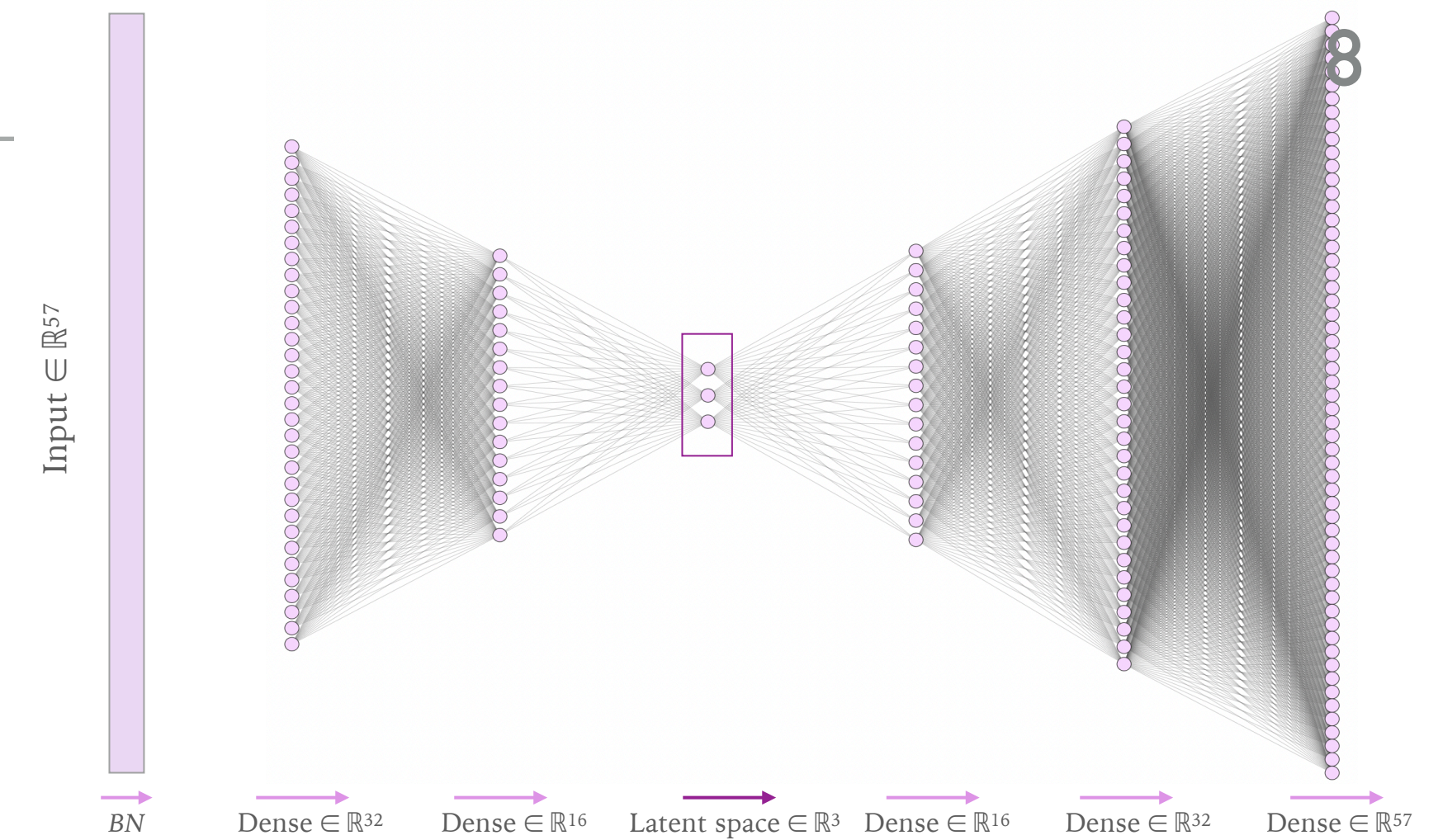BatchNorm

Output:
Conv2d 4 (1,(

ReLU    ReLU    ReLU    ReLU    ReLU

# ML IN THE TRIGGER

▸ (Variational) autoencoders for anomaly detection

▸ 1D convolutional neural networks for b-tagging



Input $\in \mathbb{R}^{57}$
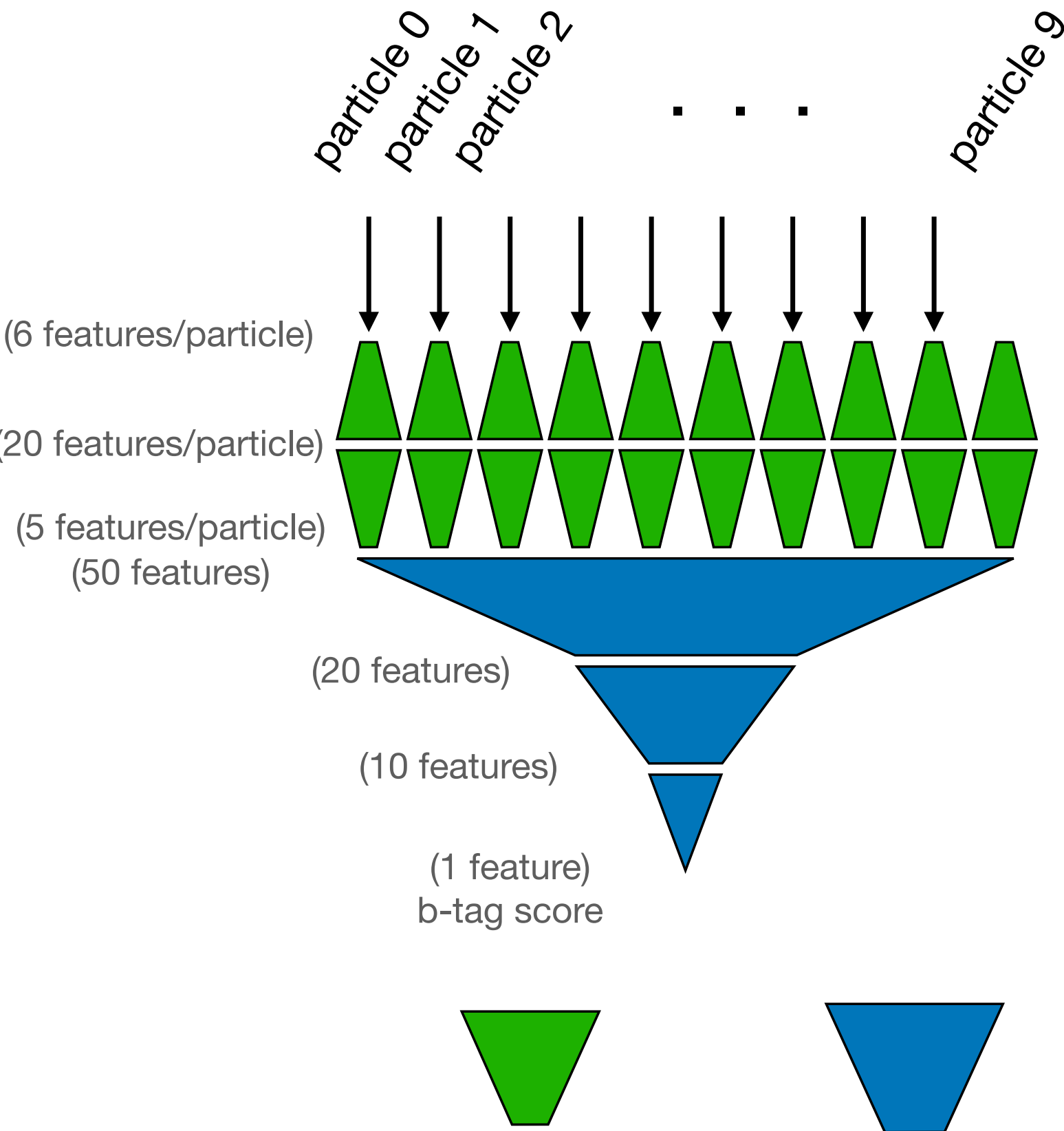
$BN$    Dense $\in \mathbb{R}^{32}$    Dense $\in \mathbb{R}^{16}$    Latent space $\in \mathbb{R}^{3}$    Dense $\in \mathbb{R}^{16}$    Dense $\in \mathbb{R}^{32}$    Dense $\in \mathbb{R}^{57}$

Block 1:
Conv2d (16,(3,3))
ReLU
AvPooling (3,1)

Block 2:
Conv2d 1 (32,(3,1))
ReLU
AvPooling (3,1)
Flatten (64)

Block 3:
Dense (8)
Dense 1 (64)
ReLU
Reshape (2,1,32)

Block 4:
Conv2d 2 (32,(3,1))
ReLU
UpSampling (3,1)
ZeroPad (0,0),(1,1)

Block 5:
Conv2d 3 (16,(3,1))
ReLU
UpSampling (3,1)
ZeroPad (1,0),(0,0)

Block 0:
Input 19x3x1
ZeroPadding (1,0)
BatchNorm

Output:
Conv2d 4 (1,(

particle 0   particle 1   particle 2   . . .   particle 9

(6 features/particle)

(20 features/particle)

(5 features/particle)
(50 features)

(20 features)

(10 features)

(1 feature)
b-tag score

Pointwise convolution
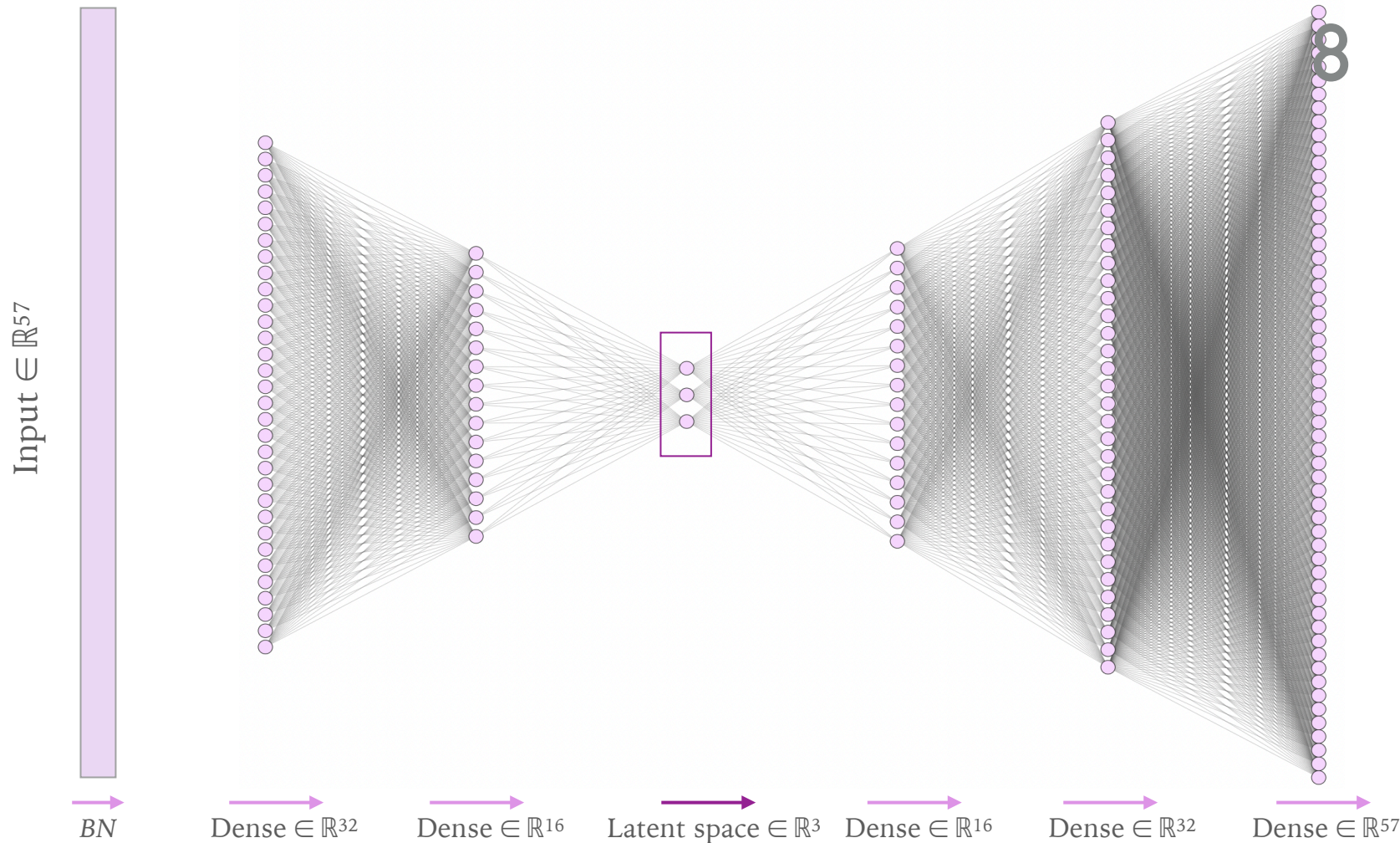(per particle dense layer)

Dense layer

# ML IN THE TRIGGER

- ▸ (Variational) autoencoders for anomaly detection
- ▸ 1D convolutional neural networks for b-tagging
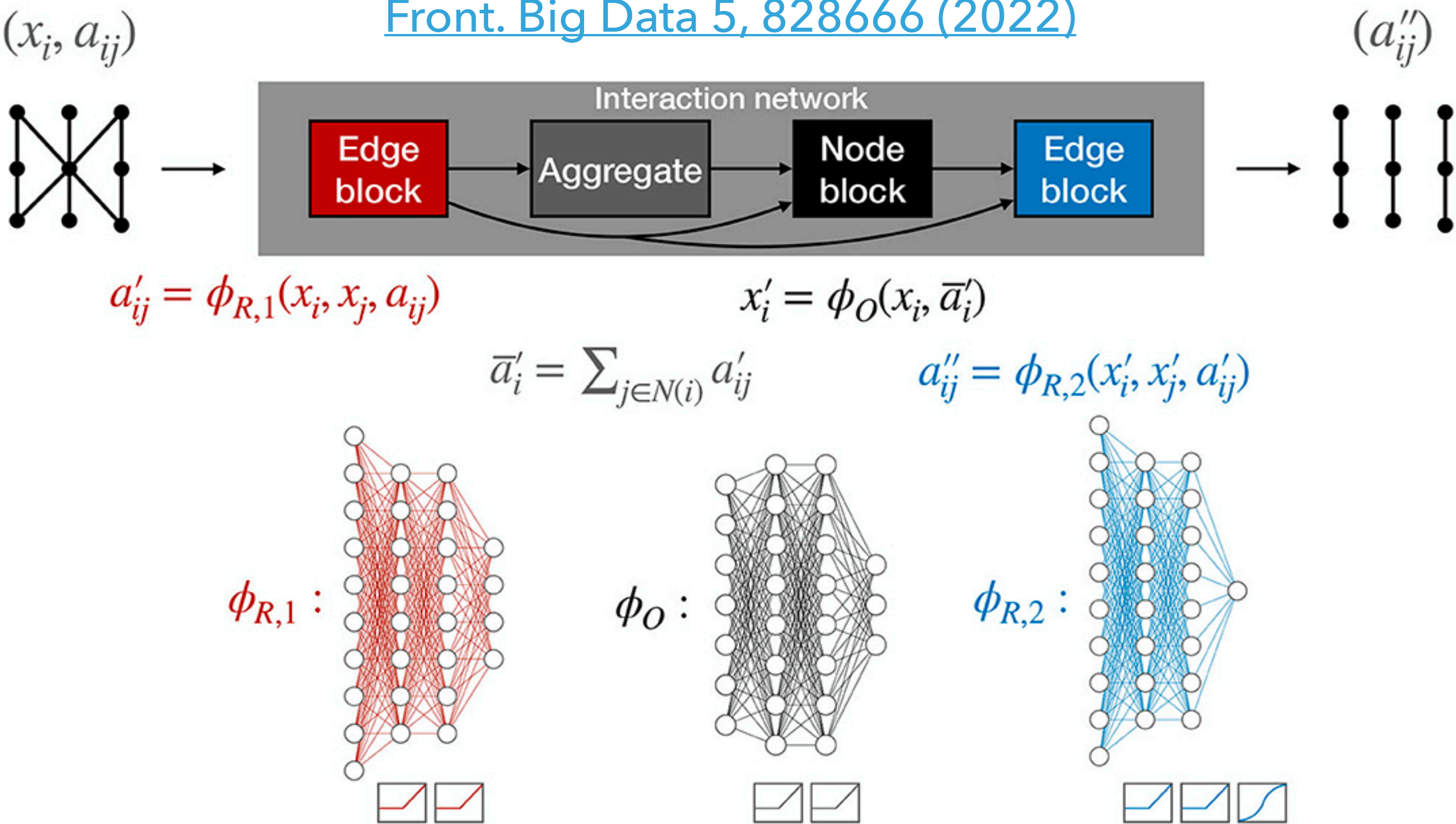- ▸ Graph neural networks for tracking



Input $\in \mathbb{R}^{57}$

$BN$   Dense $\in \mathbb{R}^{32}$   Dense $\in \mathbb{R}^{16}$   Latent space $\in \mathbb{R}^3$   Dense $\in \mathbb{R}^{16}$   Dense $\in \mathbb{R}^{32}$   Dense $\in \mathbb{R}^{57}$

particle 0   particle 1   particle 2   . . .   particle 9

(6 features/particle)

(20 features/particle)

(5 features/particle)
(50 features)

(20 features)

(10 features)

(1 feature)
b-tag score

Pointwise convolution
(per particle dense layer)

Dense layer

$(x_i, a_{ij})$

Edge block

$a'_{ij} = \phi_{R,1}(x_i, x_j, a_{ij})$

$\overline{a}'_i$

$\phi_{R,1}:$   $\phi_O:$   $\phi_{R,2}:$

Block 0:
Input 19x3x1
ZeroPadding (1,0)
BatchNorm

Block 1:
Conv2d (16,(3,3))
ReLU
AvPooling (3,1)

Block 2:
Conv2d 1 (32,(3,1))
ReLU
AvPooling (3,1)
Flatten (64)

Block 3:
Dense (8)
Dense 1 (64)
ReLU
Reshape (2,1,32)

Block 4:
Conv2d 2 (32,(3,1))
ReLU
UpSampling (3,1)
ZeroPad (0,0),(1,1)

Block 5:
Conv2d 3 (16,(3,1))
ReLU
UpSampling (3,1)
ZeroPad (1,0),(0,0)

Output:
Conv2d 4 (1,0

ReLU   ReLU   ReLU   ReLU   ReLU

Event 1 ⟶ **L1 TRIGGER ALGORITHMS** **PASS**

# WHAT MAKES THIS HARD?

▸ Reconstruct all events and reject 98% of them in ~10 µs

Event 1 ⟶ | L1 TRIGGER ALGORITHMS | PASS |

Event 2 ⟶ | L1 TRIGGER ALGORITHMS | FAIL |

Event 3 ⟶ | L1 TRIGGER ALGORITHMS | FAIL |

⋮

▸ Reconstruct all events and reject 98% of them in ~10 μs

    ▸ Algorithms have to be <1 μs and process new events every (25 ns) $\times$ $N_{tmux}$



Latency ~ 10 μs

Event 1 → L1 TRIGGER ALGORITHMS PASS

Event 2 → L1 TRIGGER ALGORITHMS FAIL

Event 3 → L1 TRIGGER ALGORITHMS FAIL

Initiation interval = 25 ns

▸ Reconstruct all events and reject 98% of them in ~10 µs

   ▸ Algorithms have to be <1 µs and process new events every (25 ns) $\times N_{tmux}$

▸ Latency necessitates all **FPGA** design

▸ Reconstruct all events and reject 98% of them in ~10 μs

   ▸ Algorithms have to be <1 μs and process new events every (25 ns) $\times$ $N_{tmux}$

▸ Latency necessitates all **FPGA** design

   ▸ Algorithms have to fit on <1 FPGA



Programmable interconnects

Latency ~ 10 μs

Event 1 → | L1 TRIGGER ALGORITHMS | PASS |

Event 2 → | L1 TRIGGER ALGORITHMS | FAIL |

Event 3 → | L1 TRIGGER ALGORITHMS | FAIL |

Initiation interval = 25 ns

▸ Reconstruct all events and reject 98% of them in ~10 μs

  ▸ Algorithms have to be <1 μs and process new events every (25 ns) × $N_{tmux}$

▸ Latency necessitates all **FPGA** design

  ▸ Algorithms have to fit on <1 FPGA

▸ How can we satisfy these constraints?



Programmable interconnects

Latency ~ 10 μs

Event 1 → L1 TRIGGER ALGORITHMS | PASS

Event 2 → L1 TRIGGER ALGORITHMS | FAIL

Event 3 → L1 TRIGGER ALGORITHMS | FAIL

Initiation interval = 25 ns

▸ **Codesign**: intrinsic development loop between ML design, training, and implementation

▸ Pruning

    ▸ Maintain high performance while removing redundant operations

▸ **Codesign**: intrinsic development loop between ML design, training, and implementation

▸ Pruning

▸ Maintain high performance while removing redundant operations

▸ Quantization

▸ Reduce precision from 32-bit floating point to 16-bit, 8-bit, …

▸ **Codesign**: intrinsic development loop between ML design, training, and implementation

▸ Pruning

　▸ Maintain high performance while removing redundant operations

▸ Quantization

　▸ Reduce precision from 32-bit floating point to 16-bit, 8-bit, …

▸ Parallelization

　▸ Balance parallelization (how fast) with resources needed (how costly)

# I. INTRO & MOTIVATION
# II. COMPRESSION
# III. HARDWARE
# IV. APPLICATIONS

Small NN benchmark correctly identifies particle "jets" 70-80% of the time

**16 inputs**

**64 nodes**

**32 nodes**

**32 nodes**

**5 outputs**

*u,d* or *s* jet

gluon jet

*W* or *Z* jet

top jet

**hls4ml**

▸ Train with **L₁ regularization** (down-weights unimportant synapses)

$$L_\lambda(\boldsymbol{w}) = L(\boldsymbol{w}) + \lambda \|\boldsymbol{w}\|_1 \qquad \|\boldsymbol{w}\|_1 = \sum_i |w_i|$$

▸ Train with **L₁ regularization** (down-weights unimportant synapses)

$$L_\lambda(\boldsymbol{w}) = L(\boldsymbol{w}) + \lambda \|\boldsymbol{w}\|_1 \qquad \|\boldsymbol{w}\|_1 = \sum_i |w_i|$$

▸ Remove **smallest** weights

▸ Train with **L₁ regularization** (down-weights unimportant synapses)

$$L_\lambda(\boldsymbol{w}) = L(\boldsymbol{w}) + \lambda \|\boldsymbol{w}\|_1 \qquad \|\boldsymbol{w}\|_1 = \sum_i |w_i|$$

▸ Remove **smallest** weights

▸ Iterate

**70% REDUCTION OF WEIGHTS WITH NO LOSS IN PERF.**

▸ Train with **L₁ regularization** (down-weights unimportant synapses)

$$L_\lambda(\boldsymbol{w}) = L(\boldsymbol{w}) + \lambda\|\boldsymbol{w}\|_1 \qquad \|\boldsymbol{w}\|_1 = \sum_i |w_i|$$

▸ Remove **smallest** weights

▸ Iterate

70% REDUCTION OF WEIGHTS WITH NO LOSS IN PERF.

# ITERATIVE MAGNITUDE-BASED PRUNING

Used in
CMS-DP-2022-020
for NN vertex finder

▸ Train with **L₁ regularization** (down-weights unimportant synapses)

$$L_\lambda(\boldsymbol{w}) = L(\boldsymbol{w}) + \lambda \|\boldsymbol{w}\|_1 \qquad \|\boldsymbol{w}\|_1 = \sum_i |w_i|$$

▸ Remove **smallest** weights

▸ Iterate

70% REDUCTION OF WEIGHTS WITH NO LOSS IN PERF.

▸ Can we compress the architecture as well?

‣ Can we compress the architecture as well?

‣ Knowledge distillation: training a small **student network** to emulate a larger **teacher model** or ensemble of networks
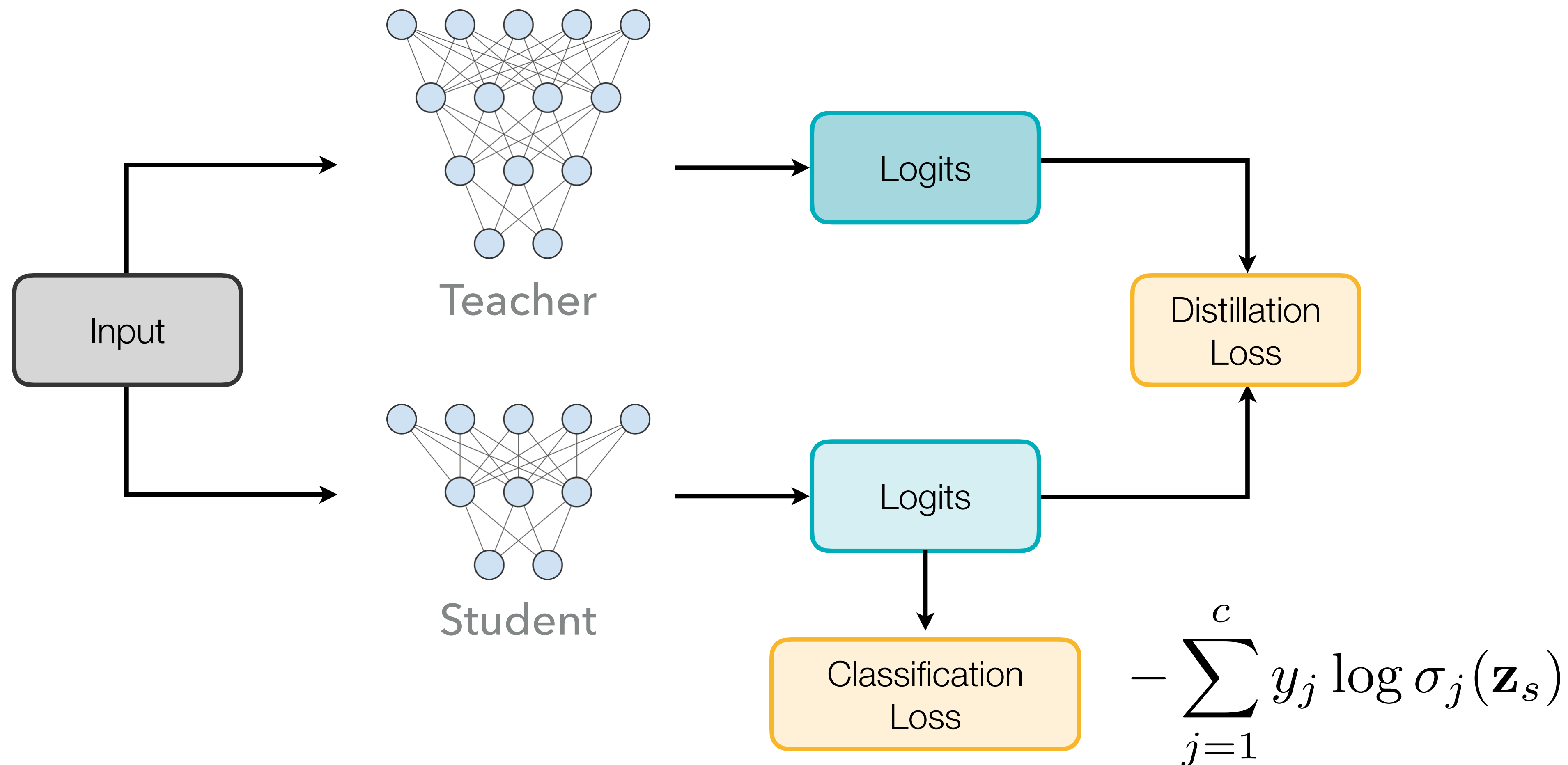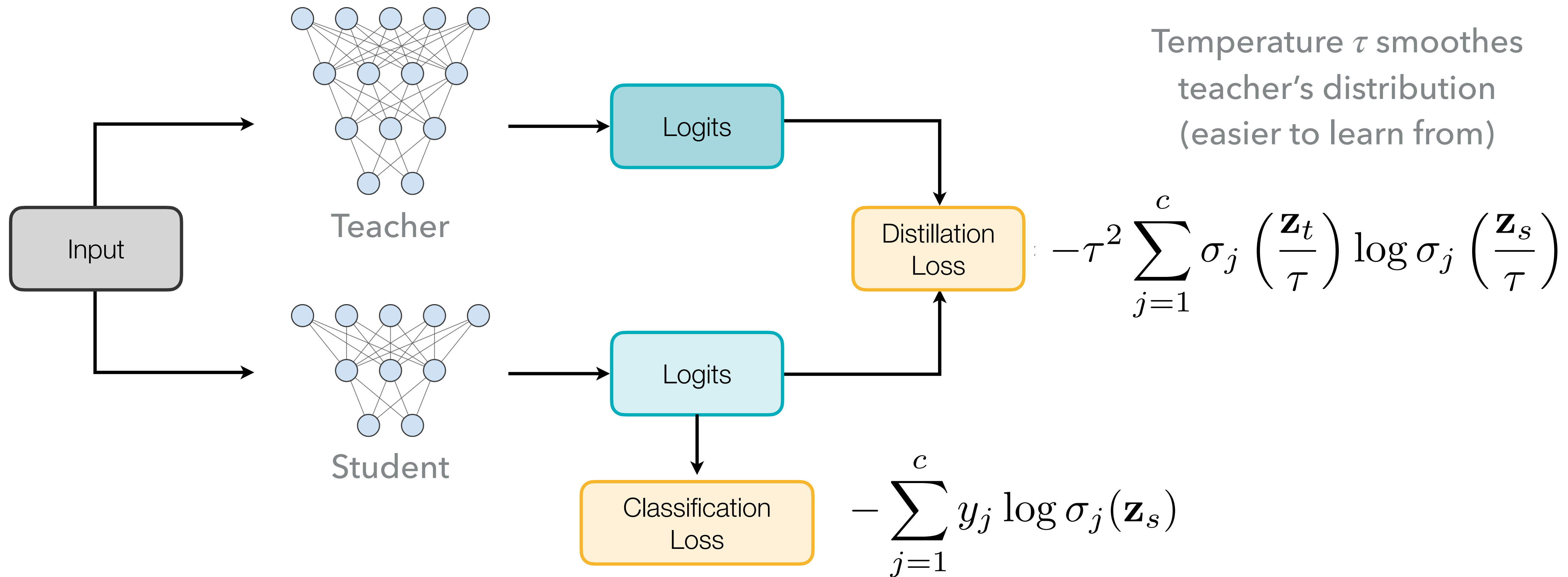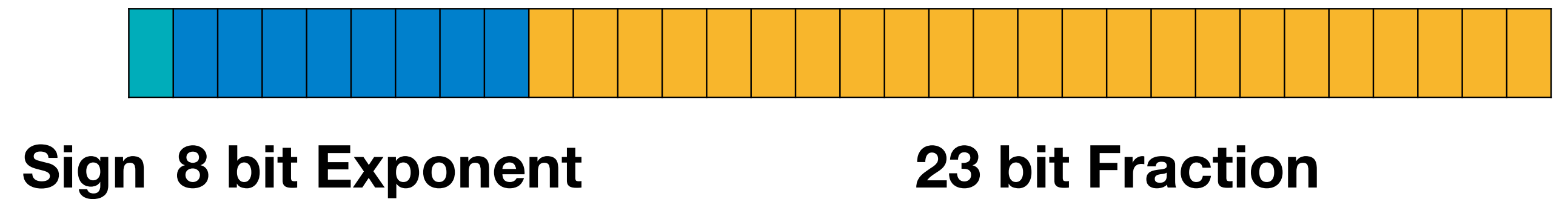
Teacher

Student

▸ Can we compress the architecture as well?

▸ Knowledge distillation: training a small **student network** to emulate a larger **teacher model** or ensemble of networks



$$-\sum_{j=1}^{c} y_j \log \sigma_j(\mathbf{z}_s)$$

‣ Can we compress the architecture as well?

‣ Knowledge distillation: training a small **student network** to emulate a larger **teacher model** or ensemble of networks



after pruning

Logits

Teacher

Input

Student

after pruning

Logits

Temperature $\tau$ smoothes teacher's distribution (easier to learn from)

Distillation Loss

$$-\tau^2 \sum_{j=1}^{c} \sigma_j \left( \frac{\mathbf{z}_t}{\tau} \right) \log \sigma_j \left( \frac{\mathbf{z}_s}{\tau} \right)$$

Classification Loss

$$-\sum_{j=1}^{c} y_j \log \sigma_j(\mathbf{z}_s)$$

‣ Quantization: using reduced precision for parameters and operations

▸ Quantization: using reduced precision for parameters and operations
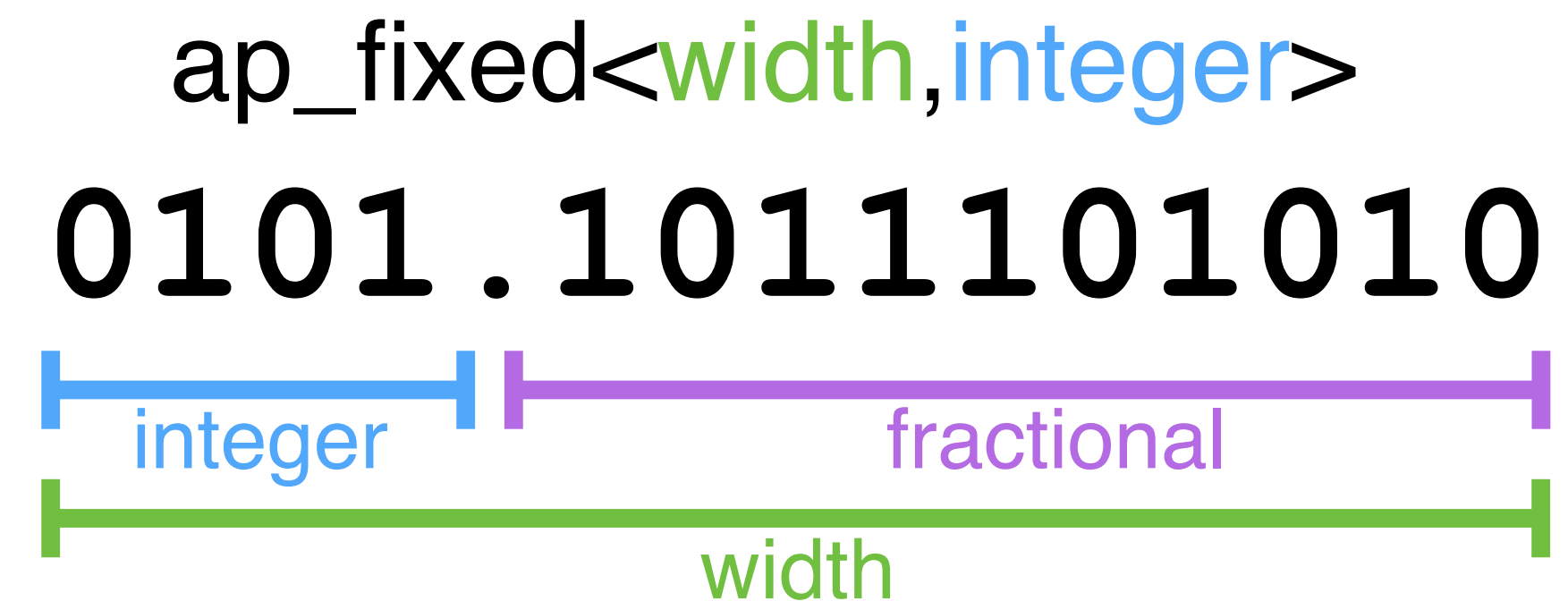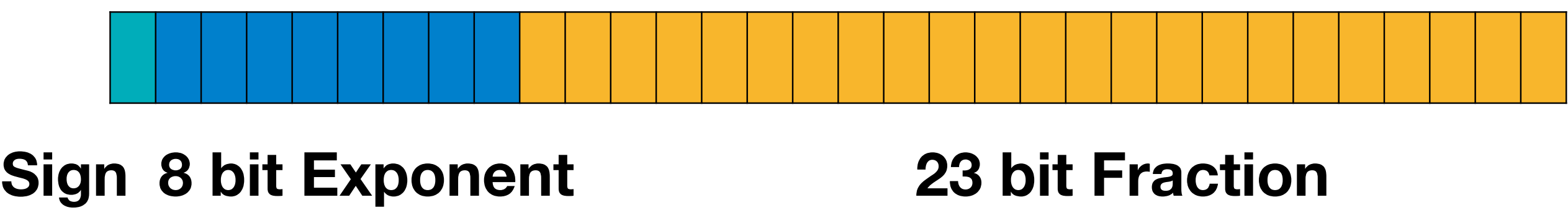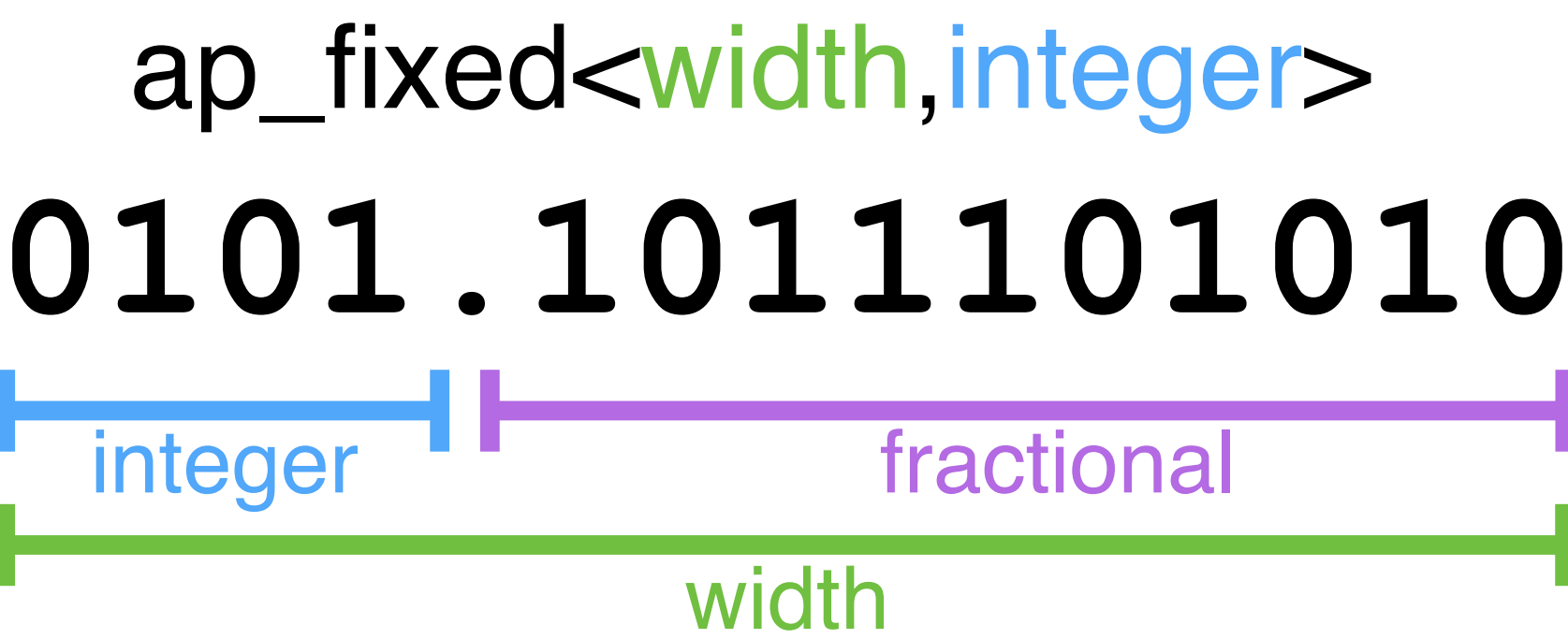
    ▸ Baseline: 32-bit floating-point precision

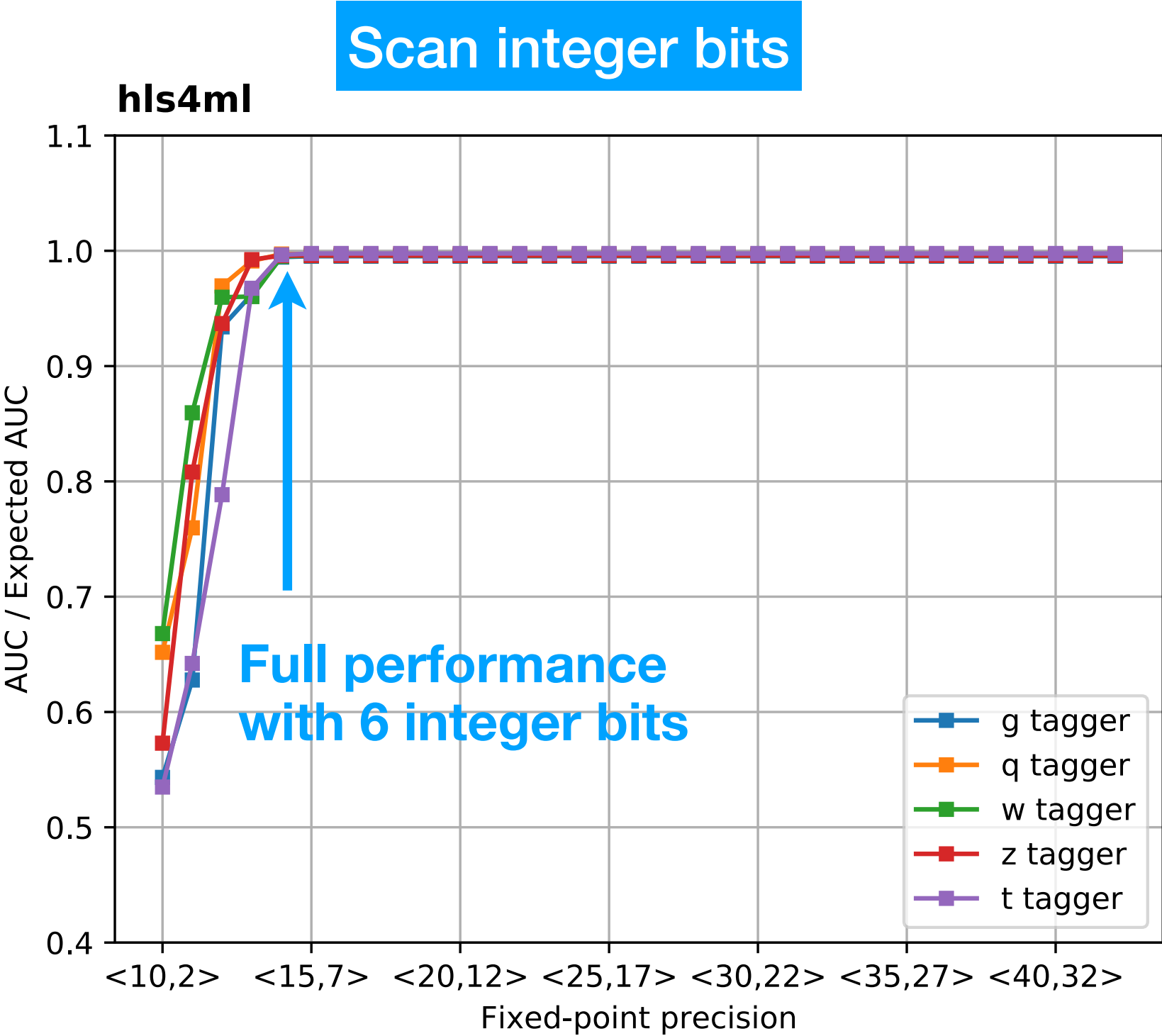**Sign  8 bit Exponent**                    **23 bit Fraction**

| 0 | 0 | 1 | 1 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

‣ Quantization: using reduced precision for parameters and operations

   ‣ Baseline: 32-bit floating-point precision

**Sign  8 bit Exponent**                    **23 bit Fraction**

‣ Fixed-point precision

ap_fixed<width,integer>

**0101.1011101010**

integer          fractional

width

0 0 1 1 1 1 1 0 1 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0

**MIT 6.S965**: TinyML and Efficient Deep

**hls4ml**

AUC

1.1

1.0

▸ Quantization: using reduced precision for parameters and operations

    ▸ Baseline: 32-bit floating-point precision

**Sign  8 bit Exponent**           **23 bit Fraction**

▸ Fixed-point precision

ap_fixed<width,integer>

**0101.1011101010**

integer          fractional

width

0 0 1 1 1 1 1 0 1 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0

$\beta = -3$        0   1        $\alpha = 4$                -$\alpha$

-128              $z$   17              127                -127

▸ Affine integer quantization

**MIT 6.S965**: TinyML and Efficient Deep

Scan integer bits

**hls4ml**



Full performance
with 6 integer bits

- g tagger
- q tagger
- w tagger
- z tagger
- t tagger

AUC / Expected AUC

Fixed-point precision

ap_fixed<width,integer>

0101.1011101010

integer   fractional
width

▸ General strategy: avoid overflows in integer bit

ap_fixed<width,integer>

0101.1011101010

integer    fractional

width

▸ General strategy: avoid overflows in integer bit

▸ Then scan the fractional bit width until reaching optimal performance

$$W^{\text{fp32}} = (S_w, W^{i4})^{\text{fp32}} \qquad a^{\text{fp32}} = W^{\text{fp32}} h^{\text{fp32}} \qquad a^{i4} = Int(\frac{a^{\text{fp32}}}{S_a})$$

Weights

INT4

INT4

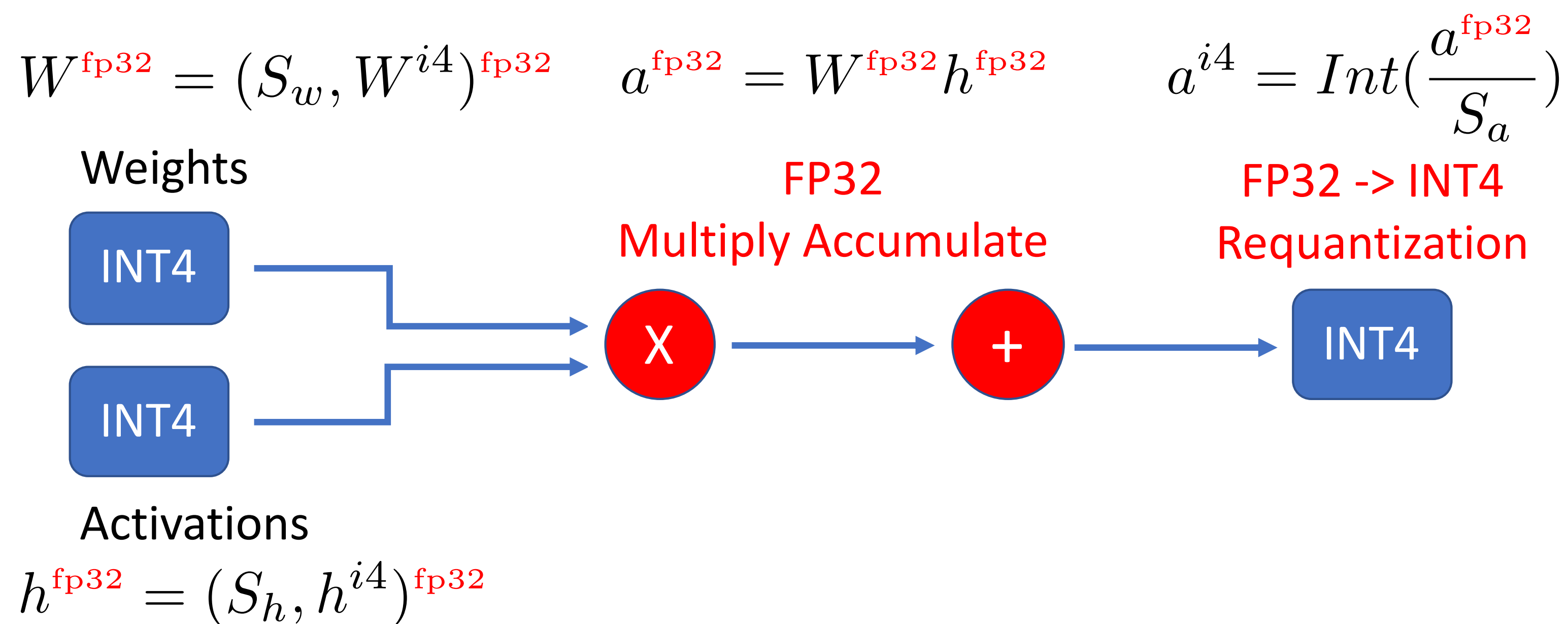Activations

$$h^{\text{fp32}} = (S_h, h^{i4})^{\text{fp32}}$$

docs.oracle.com

**What Every Computer Scientist Should Know About Floating-Point Arithmetic**

$$(S_w, W^{i4}) \qquad a^{i32} = W^{i4} h^{i4} \qquad a^{i4} = Int(\frac{S_w S_h}{S_a} a^{i32})$$

▸ Fake quantization: using 32-bit floating-point under the hood

$$W^{\mathrm{fp32}} = (S_w, W^{i4})^{\mathrm{fp32}} \qquad a^{\mathrm{fp32}} = W^{\mathrm{fp32}} h^{\mathrm{fp32}} \qquad a^{i4} = Int(\frac{a^{\mathrm{fp32}}}{S_a})$$

Weights



**docs.oracle.com**

**What Every Computer Scientist Should Know About Floating-Point Arithmetic**

Activations

$$h^{\mathrm{fp32}} = (S_h, h^{i4})^{\mathrm{fp32}}$$

$$(S_w, W^{i4}) \qquad a^{i32} = W^{i4} h^{i4} \qquad a^{i4} = Int(\frac{S_w S_h}{S_a} a^{i32})$$

▸ Fake quantization: using 32-bit floating-point under the hood

$$W^{\text{fp32}} = (S_w, W^{i4})^{\text{fp32}} \qquad a^{\text{fp32}} = W^{\text{fp32}} h^{\text{fp32}} \qquad a^{i4} = Int(\frac{a^{\text{fp32}}}{S_a})$$

Weights

INT4

FP32
Multiply Accumulate

X          +

Activations

$$h^{\text{fp32}} = (S_h, h^{i4})^{\text{fp32}}$$

docs.oracle.com

**What Every Computer Scientist Should Know About Floating-Point Arithmetic**

$$(S_w, W^{i4}) \qquad\qquad a^{i32} = W^{i4} h^{i4} \qquad a^{i4} = Int(\frac{S_w S_h}{S} a^{i32})$$

▸ Fake quantization: using 32-bit floating-point under the hood

$$W^{\text{fp32}} = (S_w, W^{i4})^{\text{fp32}} \qquad a^{\text{fp32}} = W^{\text{fp32}} h^{\text{fp32}} \qquad a^{i4} = Int(\frac{a^{\text{fp32}}}{S_a})$$

Weights

<span style="color:red">FP32
Multiply Accumulate</span>

<span style="color:red">FP32 -> INT4
Requantization</span>

**docs.oracle.com**

**What Every
Computer Scientist
Should Know About
Floating-Point
Arithmetic**

INT4

X

+

INT4

INT4

Activations

$$h^{\text{fp32}} = (S_h, h^{i4})^{\text{fp32}}$$

$$(S_w, W^{i4}) \qquad a^{i32} = W^{i4} h^{i4} \qquad a^{i4} = Int(\frac{S_w S_h}{S} a^{i32})$$

▸ Fake quantization: using 32-bit floating-point under the hood

▸ Straight-through estimator: during backpropagation, ignore quantization operation (treat as identity)

$$W^{\mathrm{fp32}} = (S_w, W^{i4})^{\mathrm{fp32}} \qquad a^{\mathrm{fp32}} = W^{\mathrm{fp32}} h^{\mathrm{fp32}} \qquad a^{i4} = Int(\frac{a^{\mathrm{fp32}}}{S_a})$$

Weights

INT4

INT4

Activations

$$h^{\mathrm{fp32}} = (S_h, h^{i4})^{\mathrm{fp32}}$$

FP32
Multiply Accumulate

FP32 -> INT4
Requantization

X

+

INT4



$$(S_w, W^{i4}) \qquad a^{i32} = W^{i4} h^{i4} \qquad a^{i4} = Int(\frac{S_w S_h}{S} a^{i32})$$

Xilinx VU9P

Xilinx VU9P

▸ Full performance with 6 bits instead of 14 bits

- Full performance with 6 bits instead of 14 bits

- Much smaller fraction of resources



Xilinx VU9P

Xilinx VU9P

- Full performance with 6 bits instead of 14 bits

- Much smaller fraction of resources

- Area & power scale quadratically with bit width



arXiv:2004.08906

▸ Quantization-aware pruning (QAP): iterative pruning can further reduce the hardware computational complexity of a quantized model

▸ After QAP, the 6-bit, 80% pruned model achieves a factor of 50 reduction in BOPs compared to the 32-bit, unpruned model

    ▸ Study using [Brevitas](Brevitas)



Bit operations (BOPs) definition:
[arXiv:1804.10969](arXiv:1804.10969)

Flat Loss Landscape

Floating Point values

...

4-bit Quantization

0          1          14         15

▸ Hessian of loss can provide additional guidance about quantization!

▸ Flat loss landscape: Lower bit width

Flat Loss Landscape

Floating Point values

4-bit Quantization

0        1        14        15

▸ Hessian of loss can provide additional guidance about quantization!

▸ Flat loss landscape: Lower bit width

▸ Sharp loss landscape: Higher bit width



Flat Loss Landscape

Floating Point values

4-bit Quantization

0    1    14    15

Sharp Loss Landscape

Floating Point values

8-bit Quantization

0    1    254    255

▸ Say you want to program an "adder" function on an FPGA

```
module adder(
    input  wire [4:0] a,
    input  wire [4:0] b,
    output wire [4:0] y
);

    assign y = a + b;


endmodule
```



▸ Register transfer-level (RTL)
  code is "synthesized" into gates

‣ Say you want to program an "adder" function on an FPGA

```verilog
module adder(
    input  wire [4:0] a,
    input  wire [4:0] b,
    output wire [4:0] y
);

    assign y = a + b;

endmodule
```



Adder

a

b

y

Synthesis

‣ Register transfer-level (RTL) code is "synthesized" into gates

▸ What if instead we specify an AI model



High-Level Synthesis

▶ hls4ml for scie

Model

Machi
optimi

hls4ml

hls 4 ml

HLS 4 ML

▸ hls4ml for scie

hls4ml

Keras
TensorFlow
PyTorch
...

hls 4 ml

FPGA flow

Model

Compressed
model

HLS
conversion

HLS
project

Machi
optimi

ASIC flow

Tune configuration
latency, throughput,
power, resource usage

▶ hls4ml for scie

**hls4ml**

Keras
TensorFlow
PyTorch
...

**Model**

**Compressed model**

AMD
XILINX

**FPGA flow**

intel

**HLS conversion**

**HLS project**

**Machi optimi**

**Tune configuration**
latency, throughput,
power, resource usage

**ASIC flow**

Mentor
A Siemens Business

▸ hls4ml for scie

hls4ml

hls4ml

Keras
TensorFlow
PyTorch
...

Model

Compressed
model

HLS
conversion

HLS
project

Machi
optimi

Tune configuration
latency, throughput,
power, resource usage

FPGA flow

ASIC flow

hls4ml v0.7.0
coming this week!

▸ FINN (NNs): https://finn.readthedocs.io/en/latest/

▸ Confier (BDTs): https://github.com/thesps/conifer

▸ fwXMachina (BDTs): http://fwx.pitt.edu/

▸ FlowGNN: https://github.com/sharc-lab/flowgnn



**(b) FlowGNN** Architecture with **Multiple** Node Transformation, **Multiple** Message Passing, and **parallelized** Edge Embedding
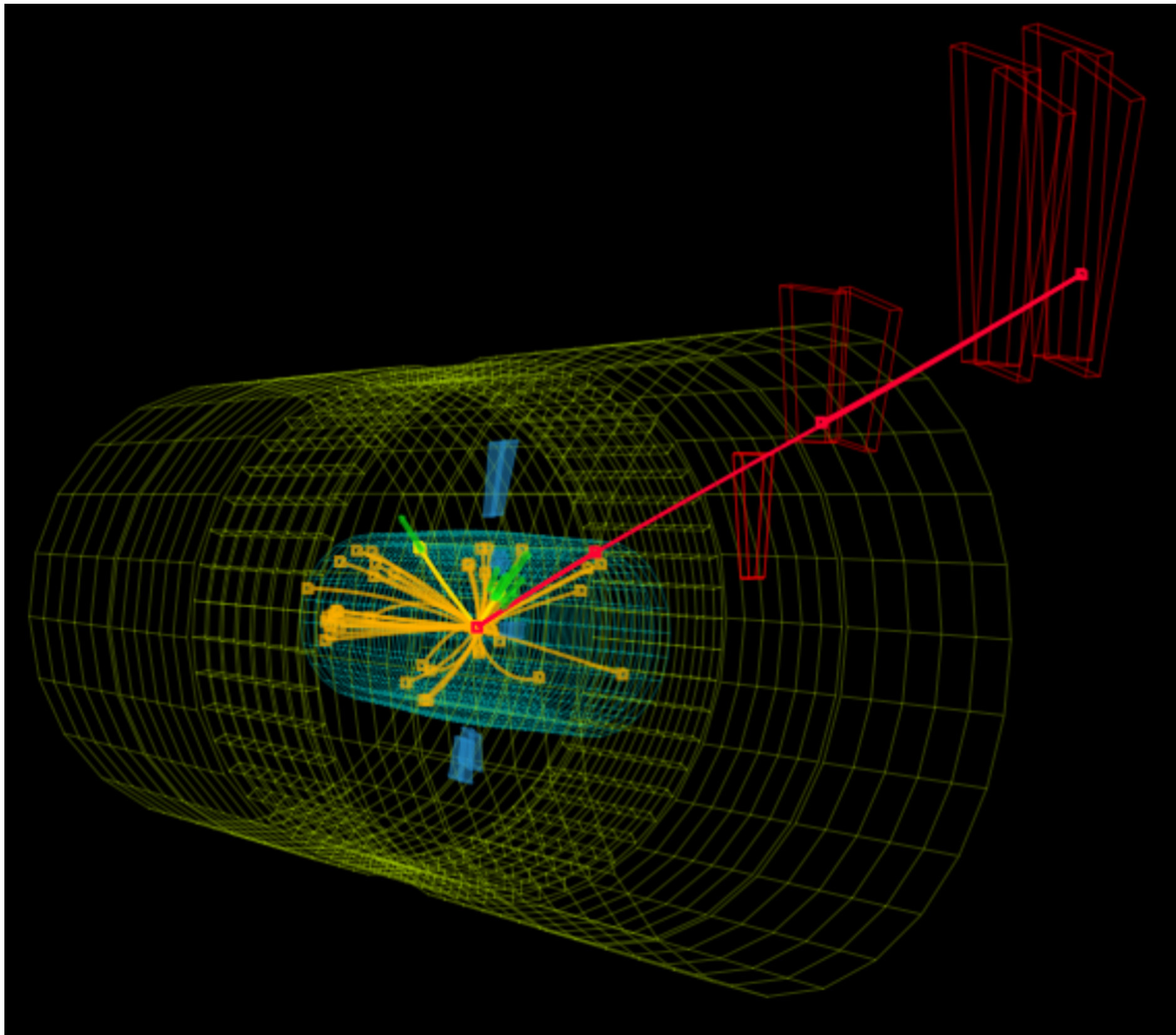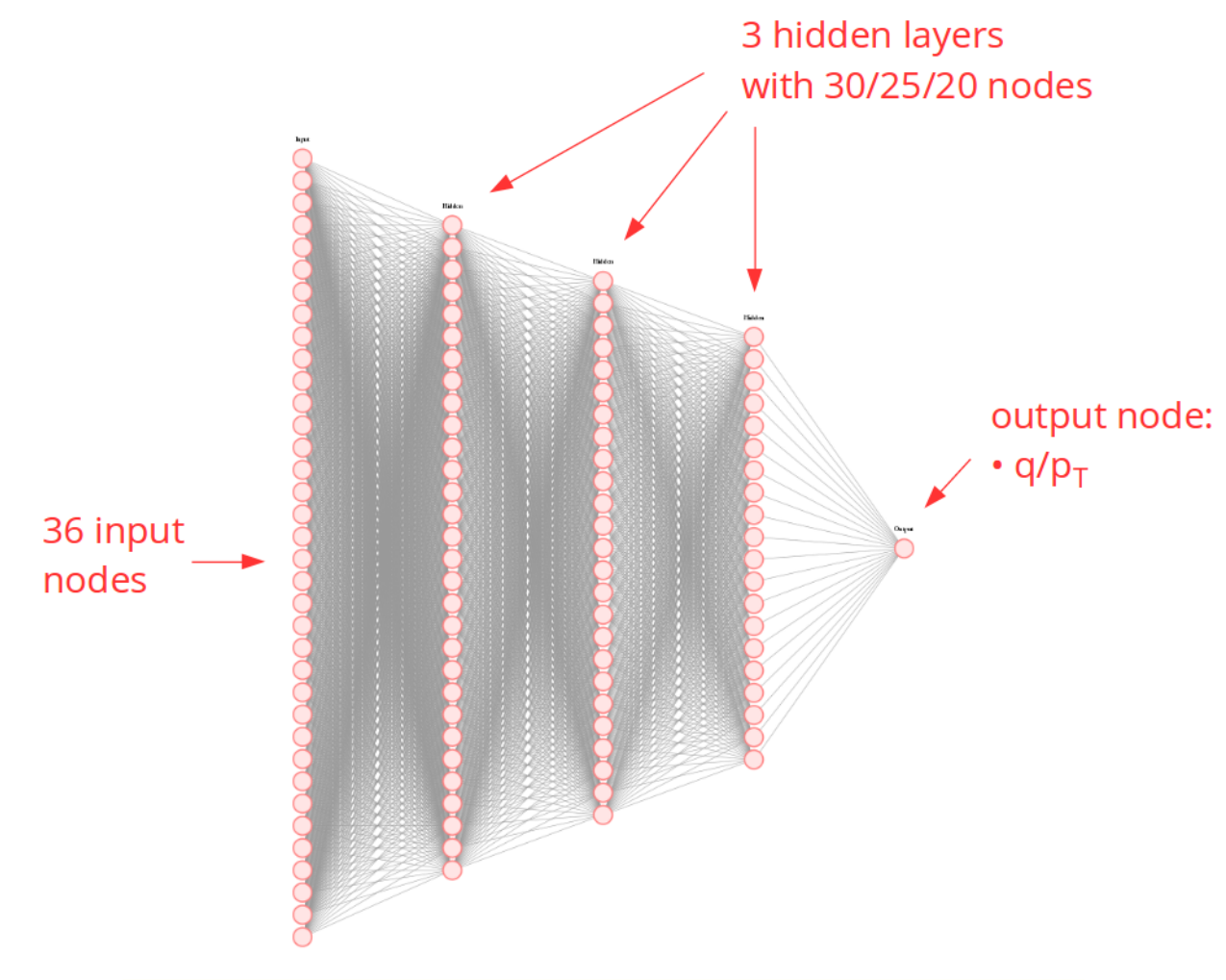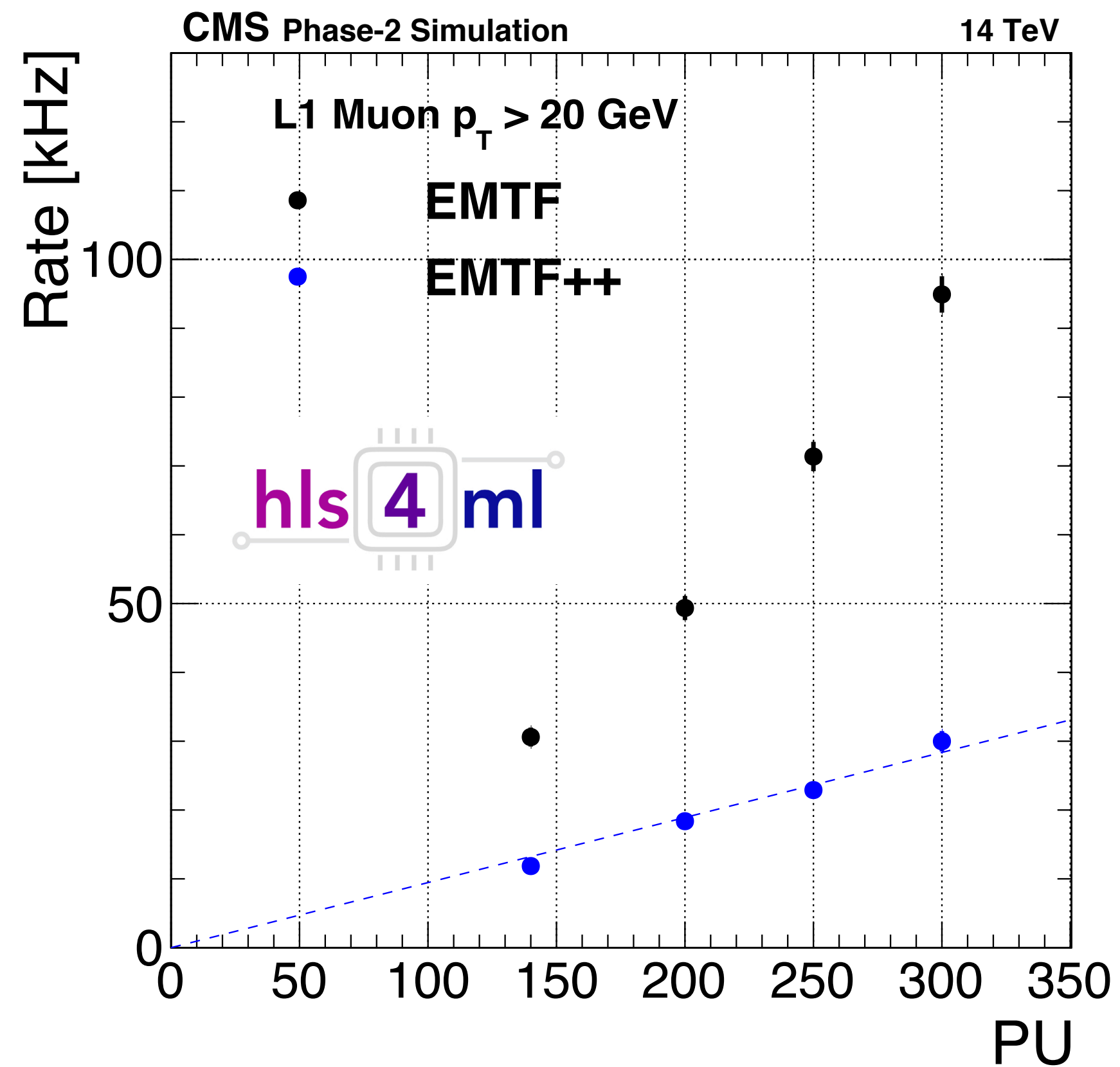
I. INTRO & MOTIVATION
II. COMPRESSION
III. HARDWARE
IV. APPLICATIONS

3 hidden layers
with 30/25/20 nodes

output node:
• q/p$_T$

36 input
nodes

3 hidden layers
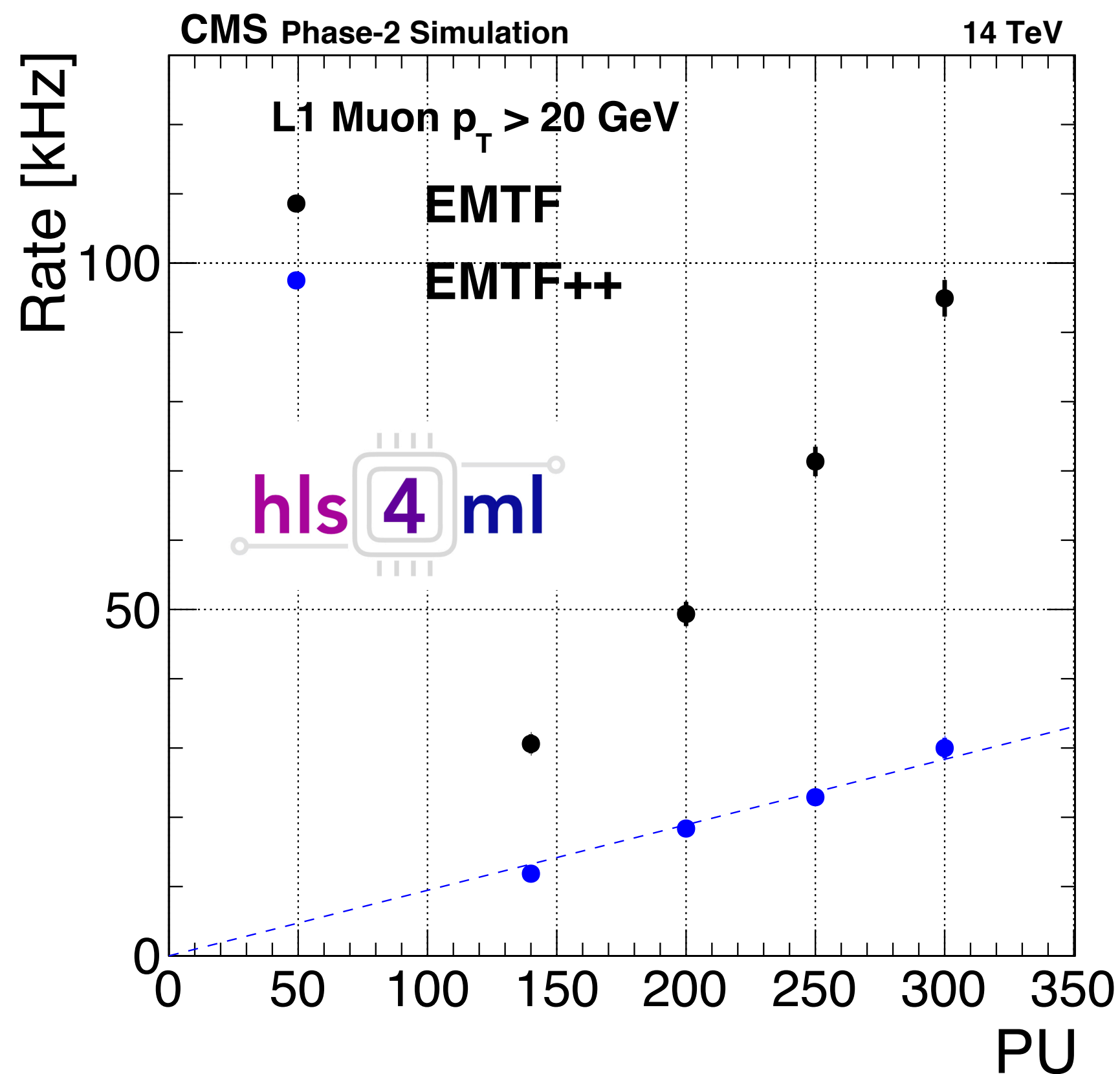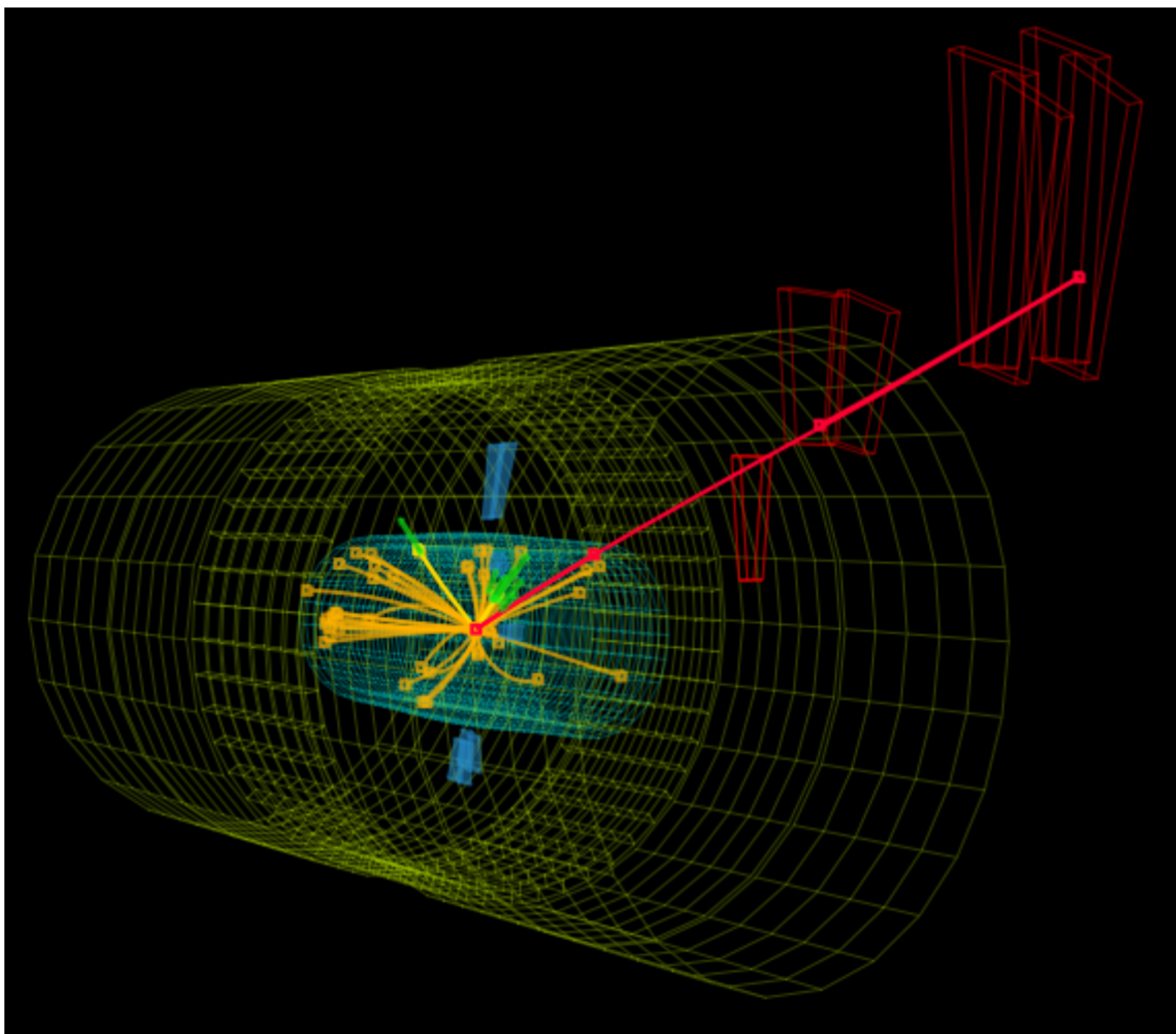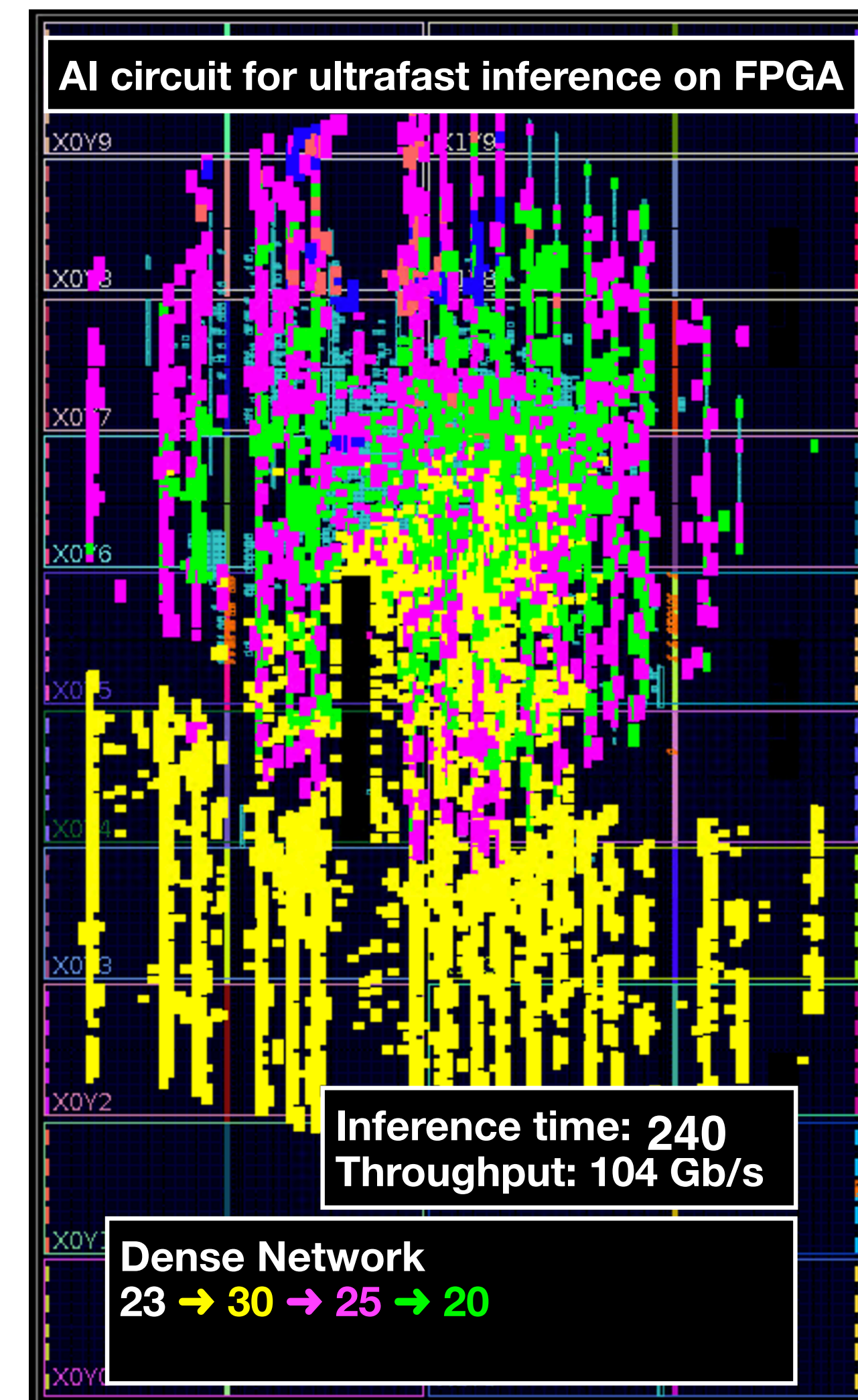with 30/25/20 nodes

36 input
nodes

output node:
• q/p$_T$

▶ NN measures muon momentum
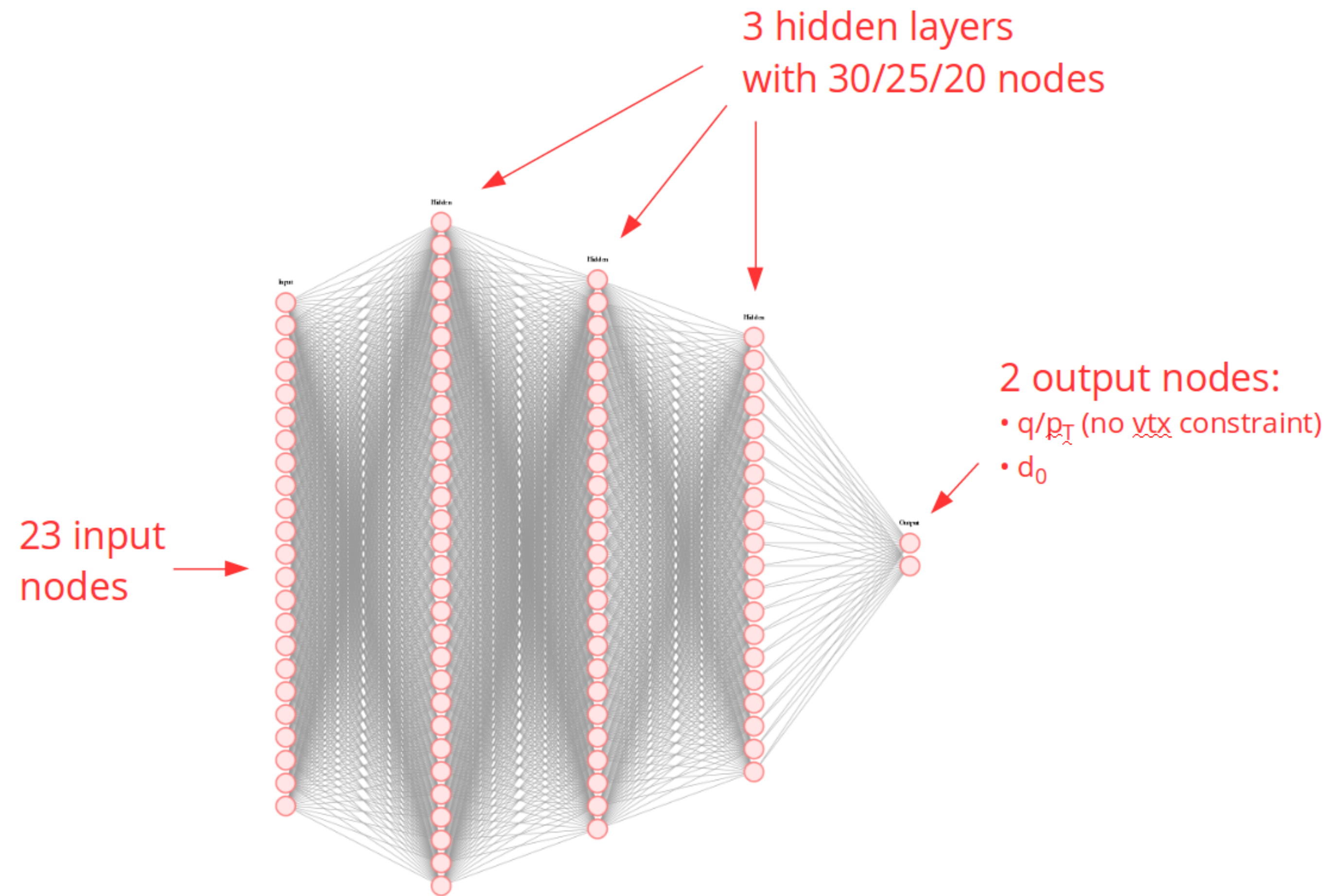
▸ NN measures muon momentum
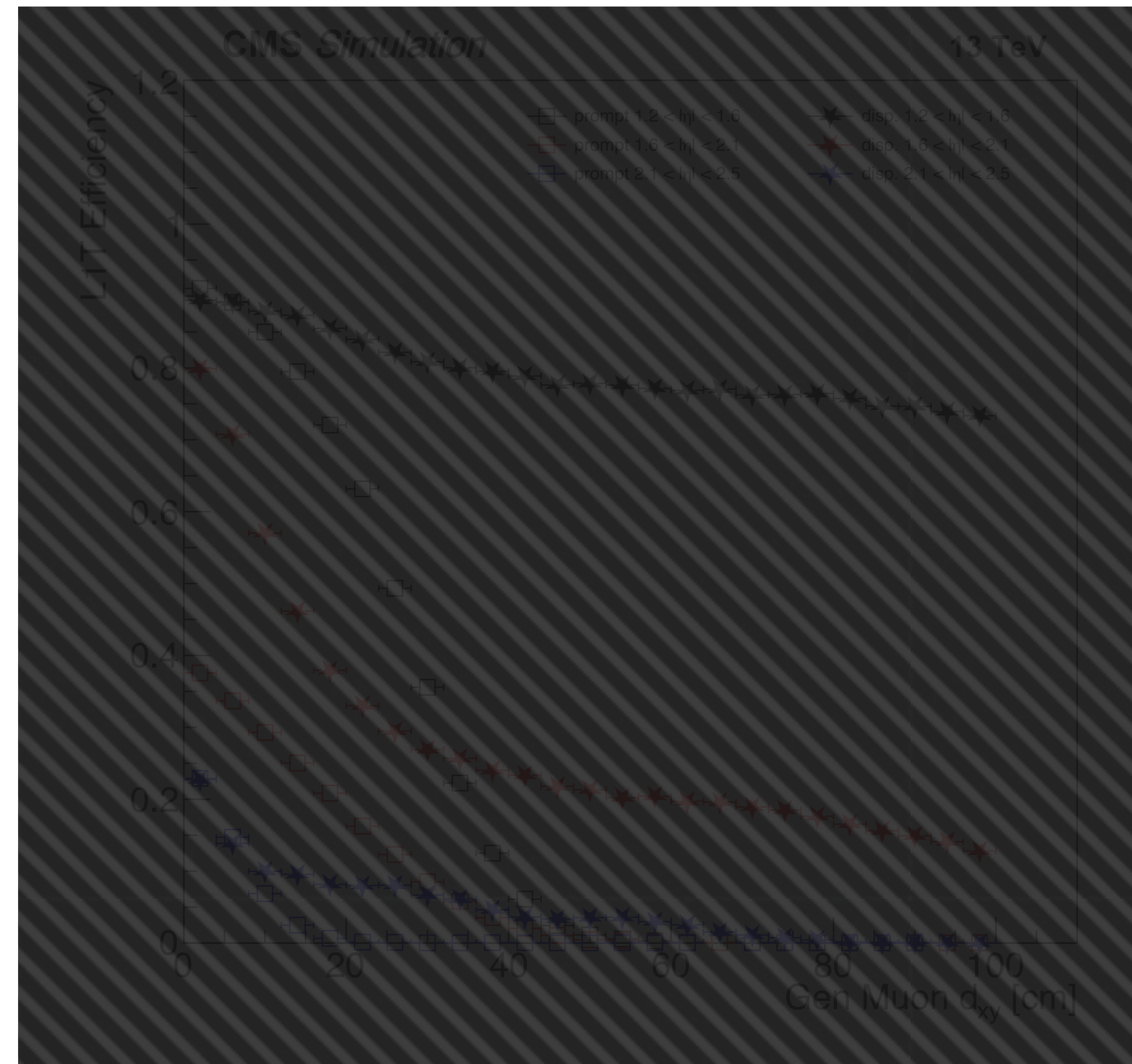
  ▸ 3× reduction in the trigger rate for NN!

▸ NN measures muon momentum

  ▸ 3× reduction in the trigger rate for NN!

▸ Fits within L1 trigger latency (240 ns!) and
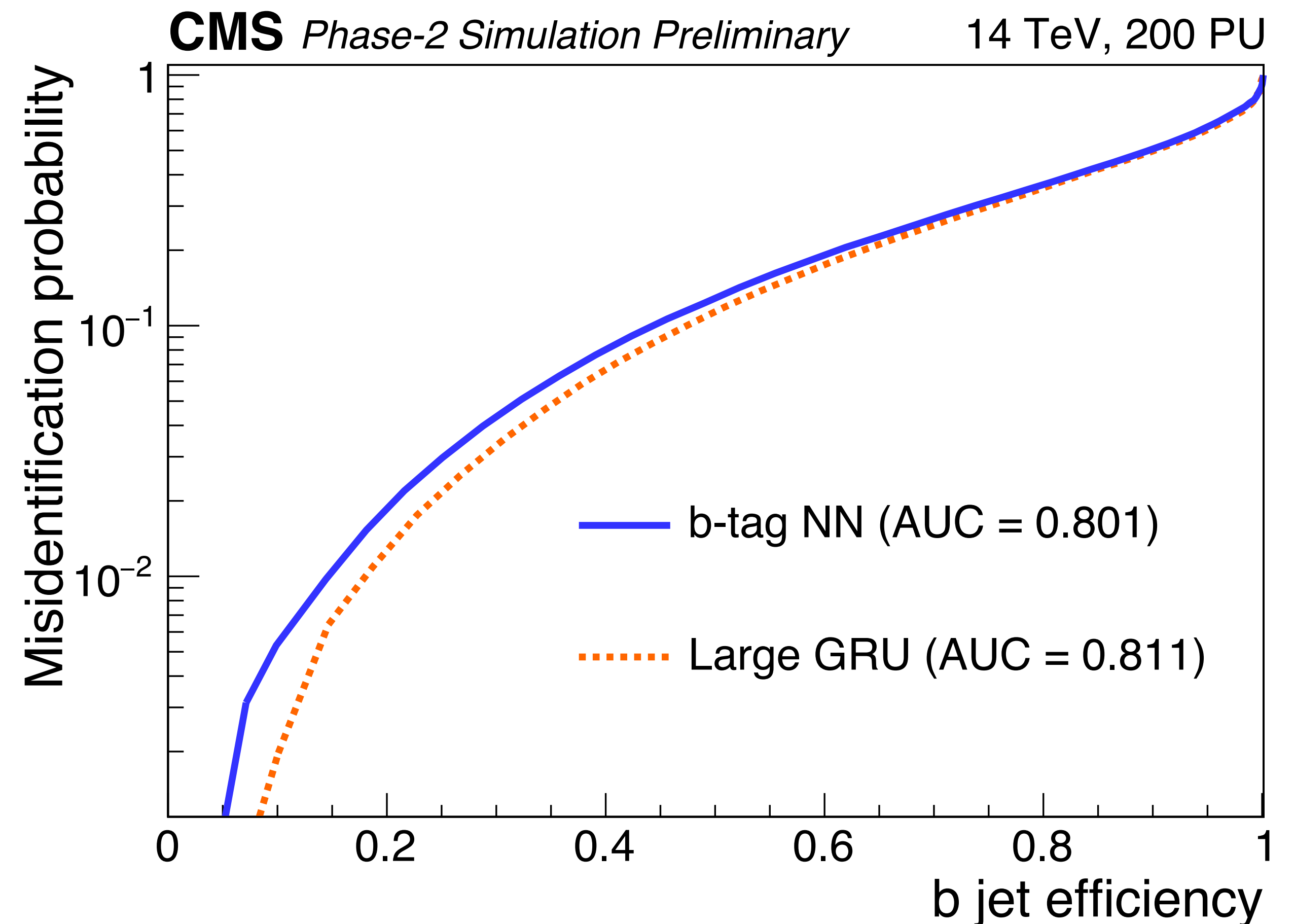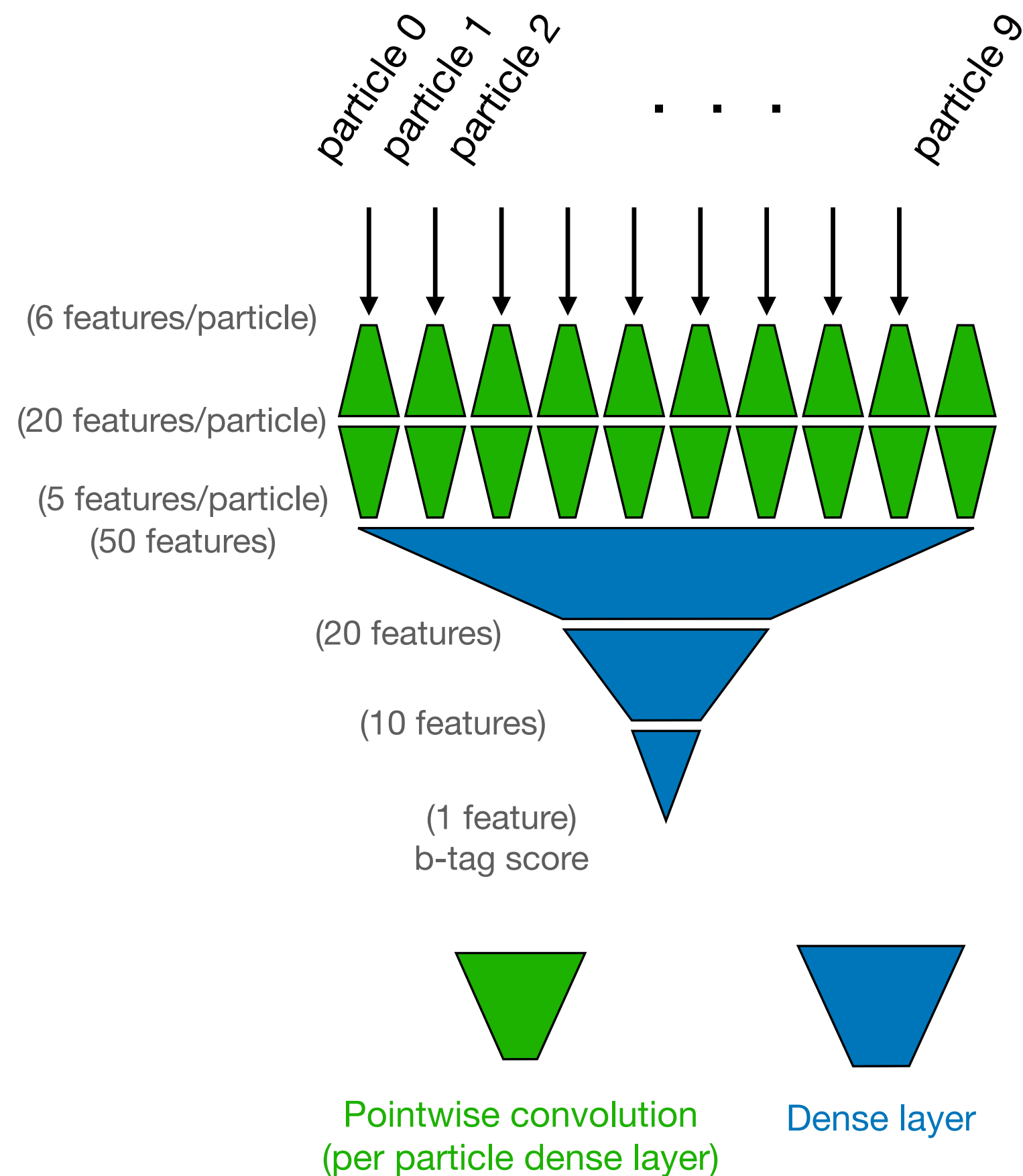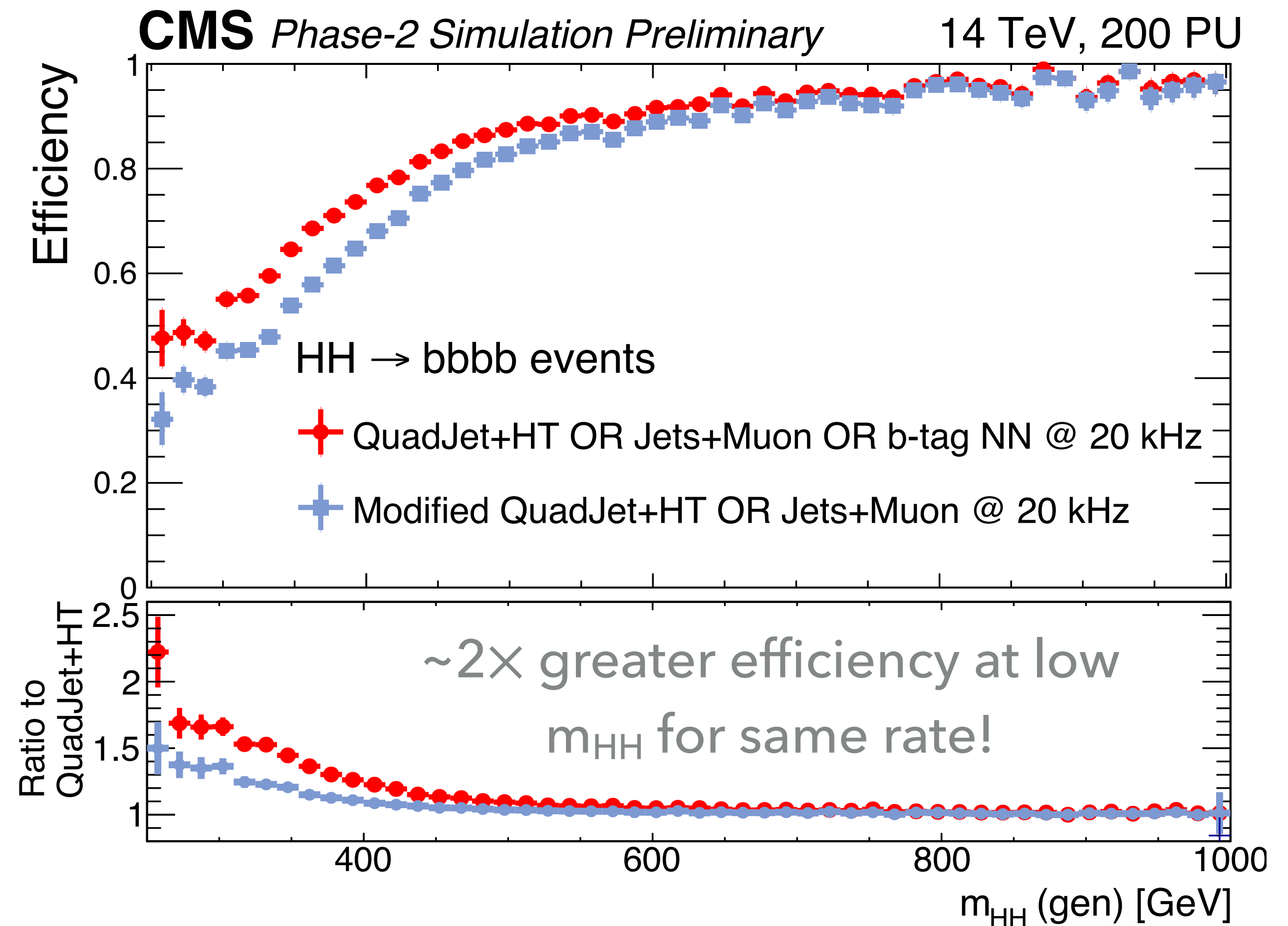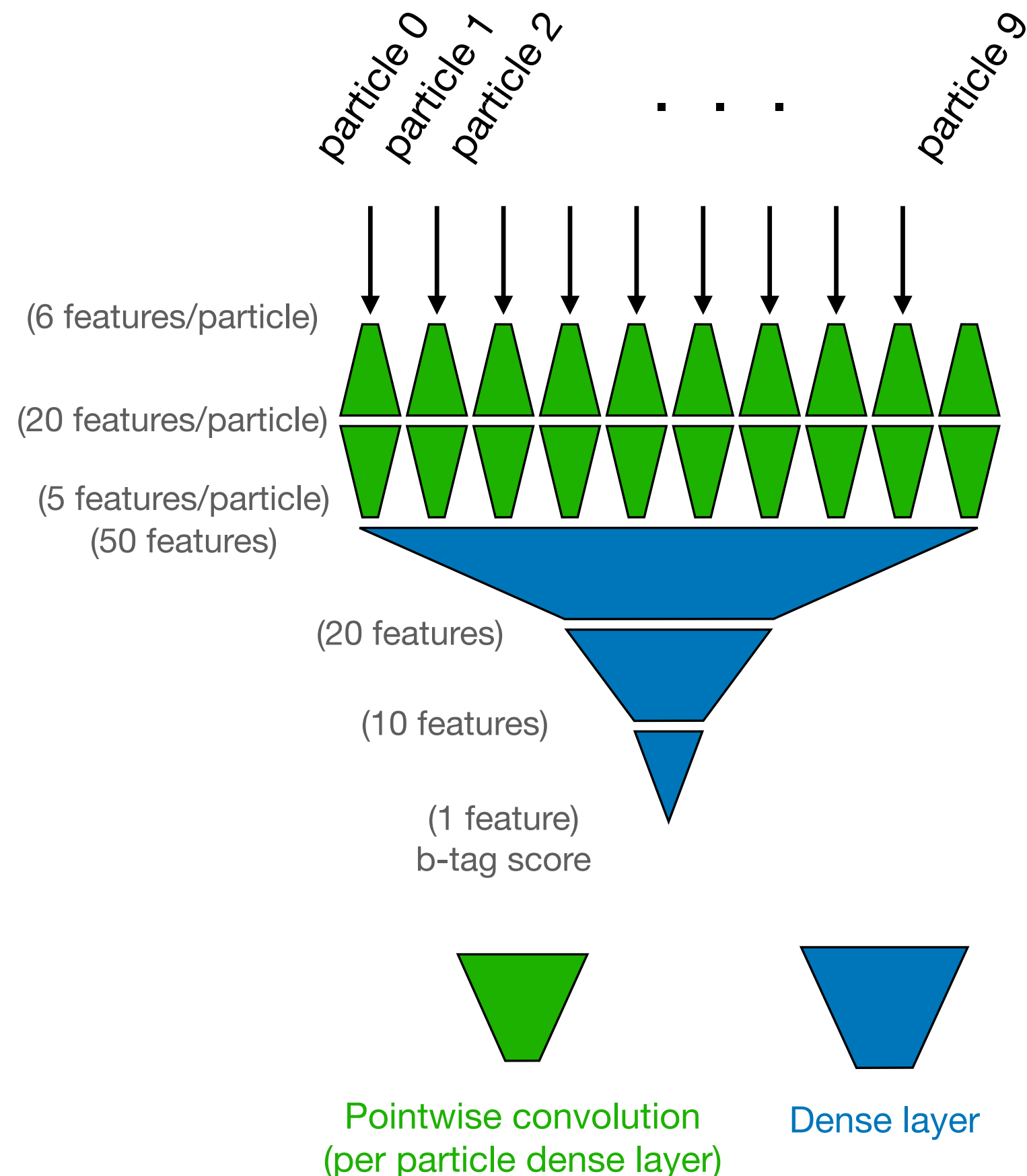  FPGA resource requirements (less then 30%)

3 hidden layers
with 30/25/20 nodes

2 output nodes:
• q/p$_T$ (no vtx constraint)
• d$_0$

23 input
nodes

displacement as well as p$_T$

▸ *Stay tuned for Run 3 results*

▶ Upgraded HL-LHC level-1 track trigger information enables b-tagging with a **neural network** to improve the $HH \to 4b$ search

  ▶ Input features for 10 particles within each jet: particle type, momentum, and vertex information



particle 0  particle 1  particle 2  . . .  particle 9

(6 features/particle)
(20 features/particle)
(5 features/particle)
(50 features)
(20 features)
(10 features)
(1 feature)
b-tag score

Pointwise convolution
(per particle dense layer)

Dense layer

**CMS** *Phase-2 Simulation Preliminary*    14 TeV, 200 PU

Misidentification probability

b jet efficiency

b-tag NN (AUC = 0.801)

Large GRU (AUC = 0.811)
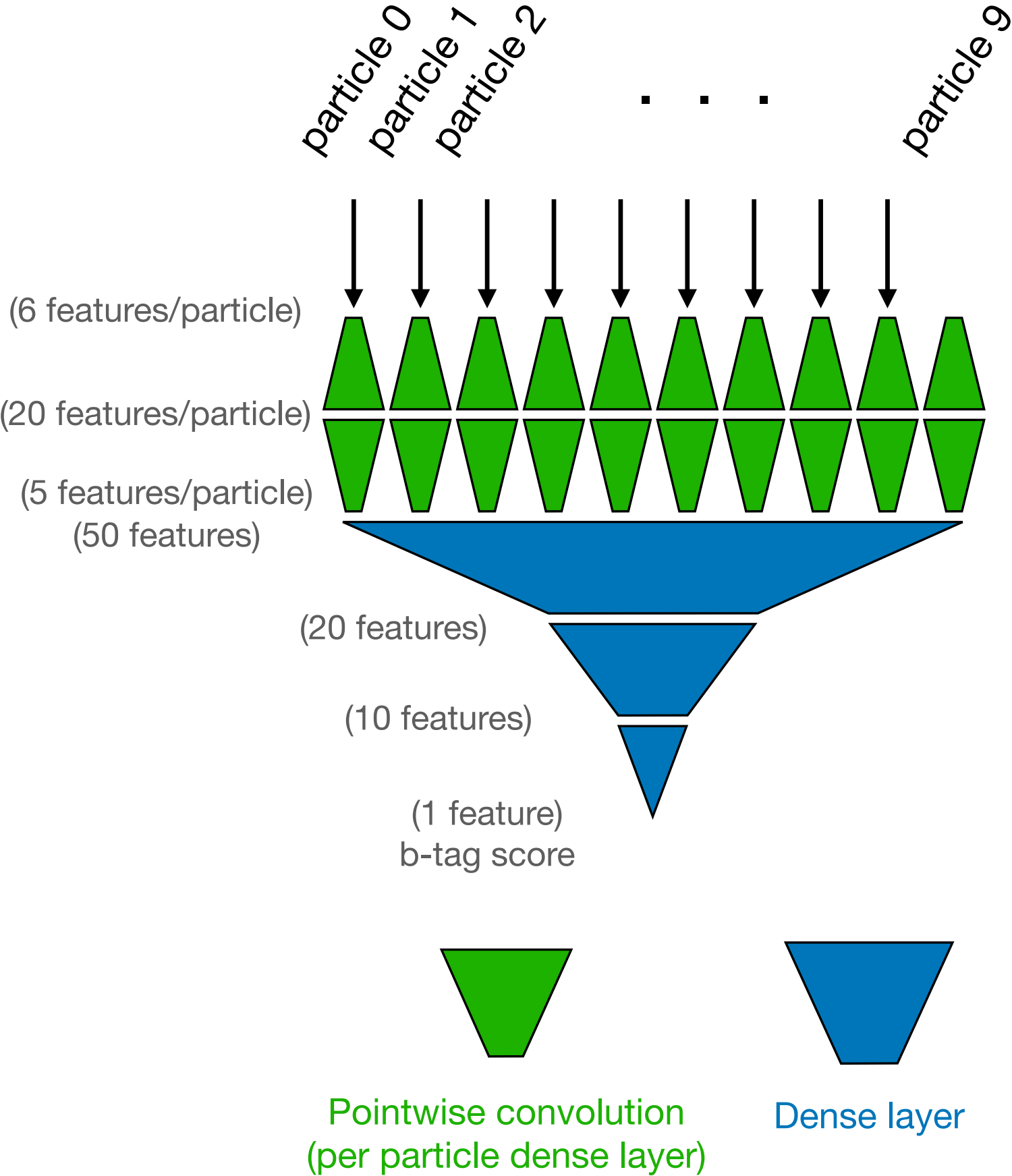
▸ Upgraded HL-LHC level-1 track trigger information enables b-tagging with a **neural network** to improve the $HH \to 4b$ search

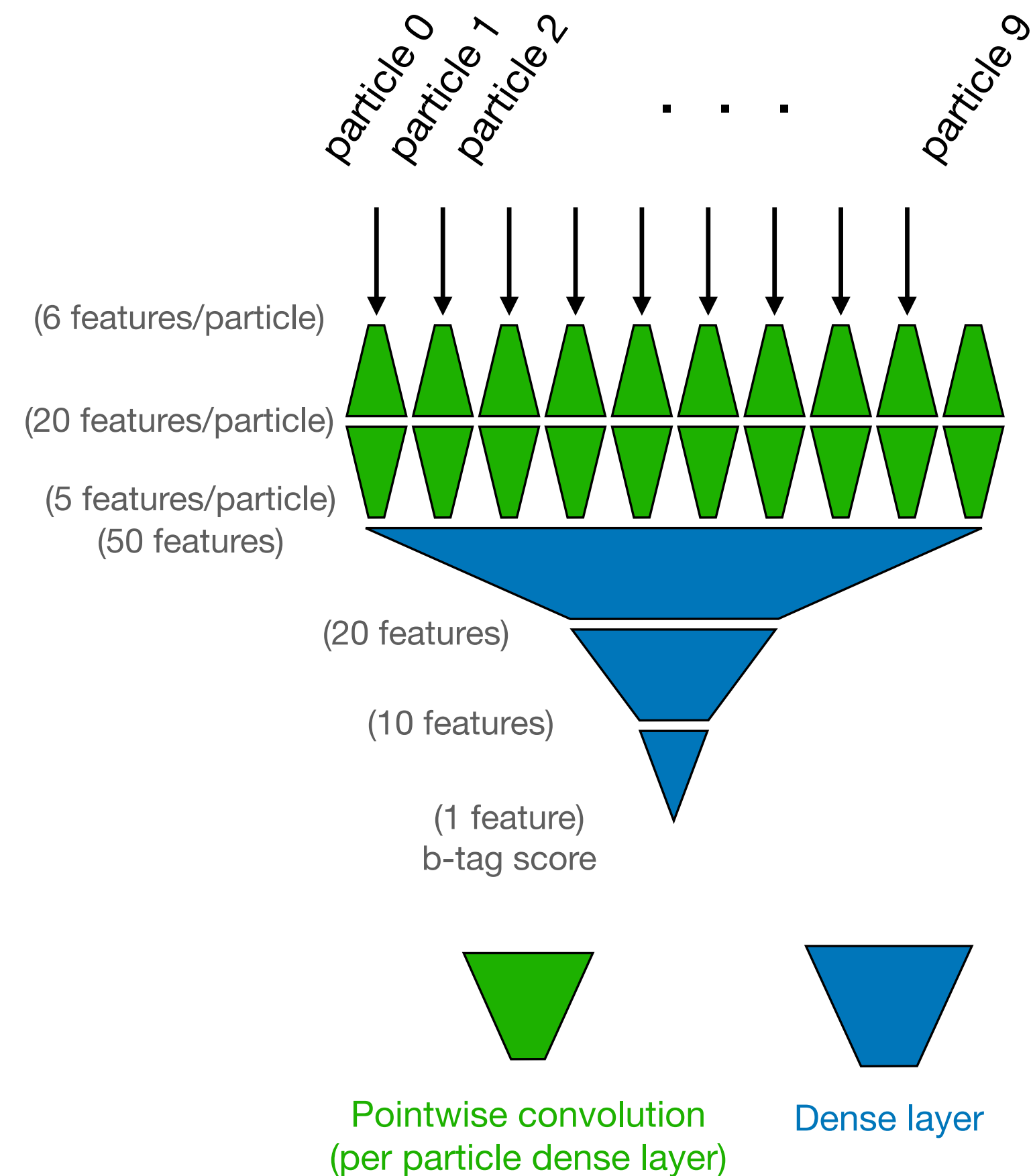   ▸ Input features for 10 particles within each jet: particle type, momentum, and vertex information

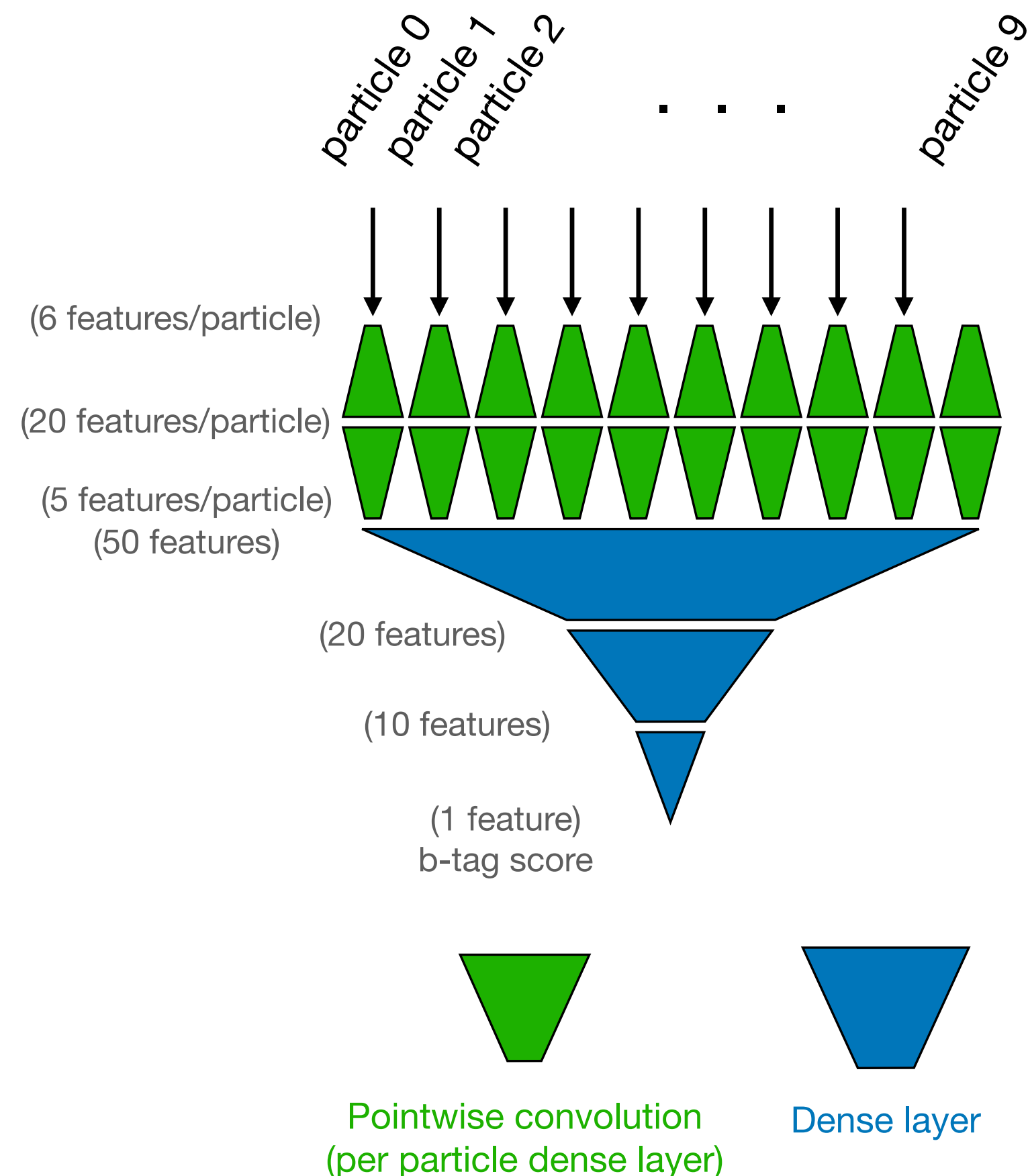(6 features/particle)

(20 features/particle)

(5 features/particle)
(50 features)

(20 features)

(10 features)

(1 feature)
b-tag score

particle 0 particle 1 particle 2 . . . particle 9

Pointwise convolution
(per particle dense layer)

Dense layer

▸ But does it fit and meet timing?



(6 features/particle)

(20 features/particle)

(5 features/particle)
(50 features)

(20 features)

(10 features)

(1 feature)
b-tag score

Pointwise convolution
(per particle dense layer)

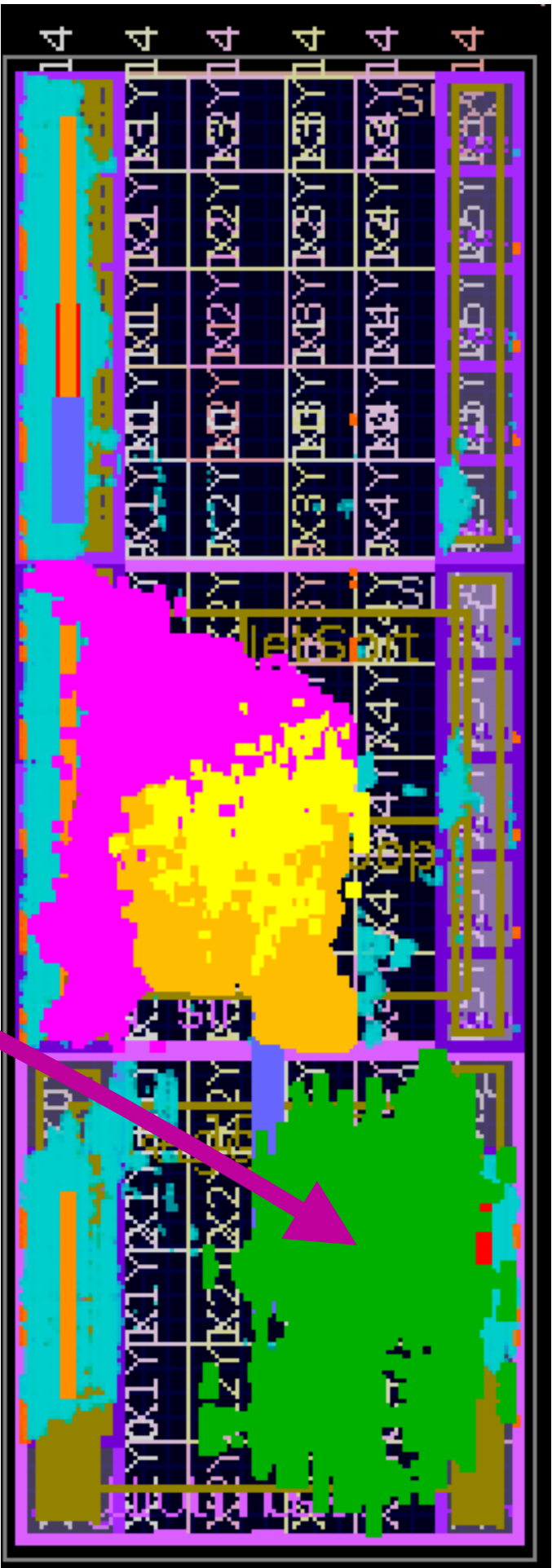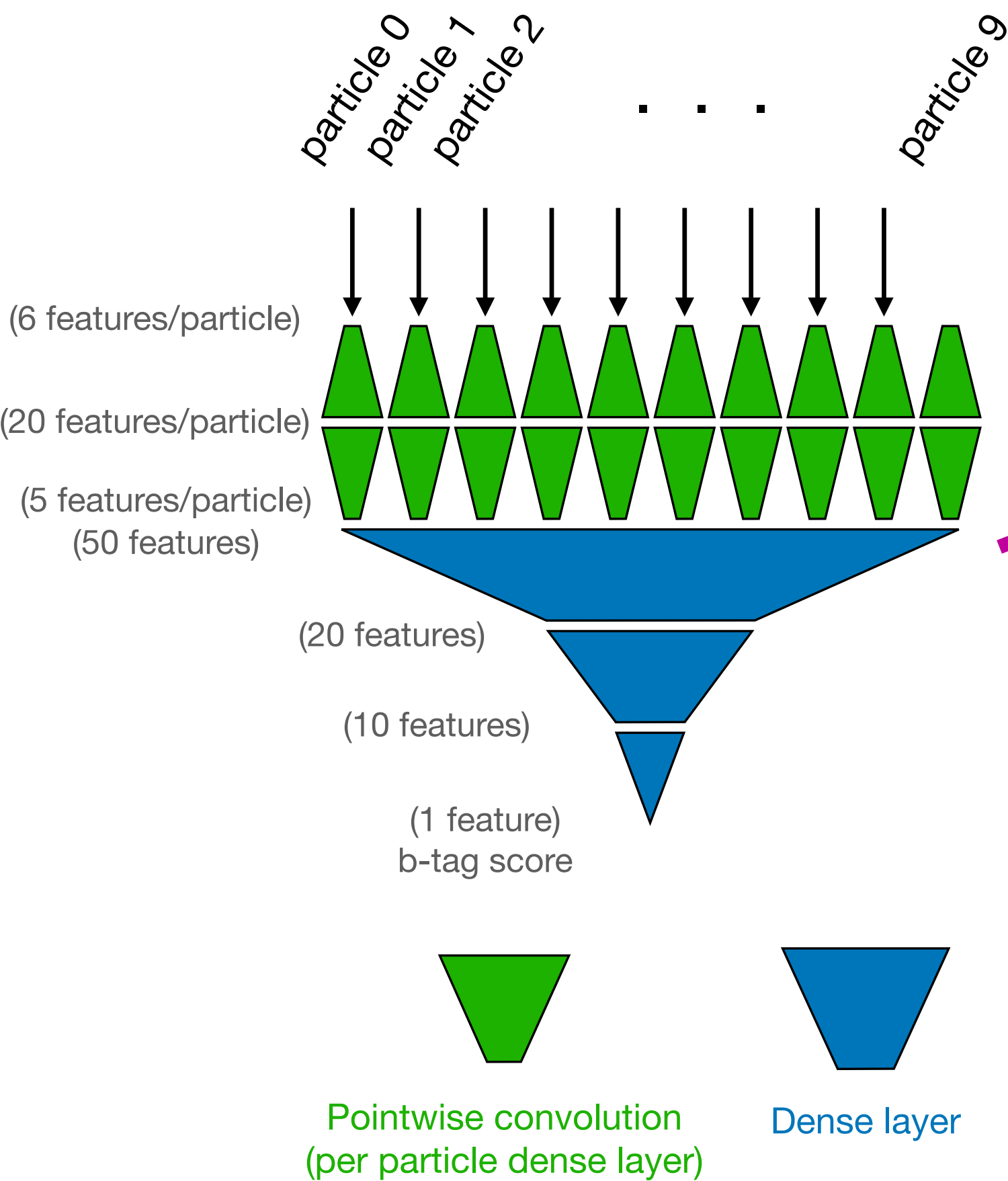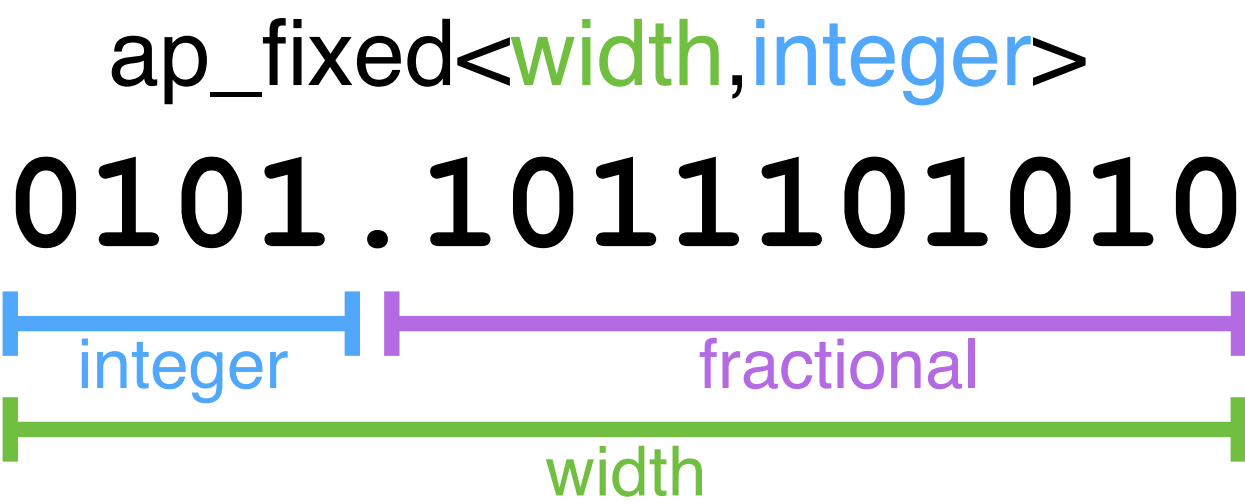Dense layer

▸ But does it fit and meet timing?

▸ After **quantization**, can implement NN with 9 bits

ap_fixed<width,integer>

**0101.1011101010**

integer | fractional

width



(6 features/particle)

(20 features/particle)

(5 features/particle)
(50 features)

(20 features)

(10 features)

(1 feature)
b-tag score

Pointwise convolution
(per particle dense layer)

Dense layer

particle 0  particle 1  particle 2  . . .  particle 9

**hls4ml**

FPGA AUC / Expected AUC

<10,2>  <15,7>  <20,12>  <25,17>  <30,22>  <35,27>

Fixed-point precision

▸ But does it fit and meet timing?

▸ After **quantization**, can implement NN with 9 bits

▸ Latency of 60 ns, II of 5 ns per jet, and <12% of FPGA

ap_fixed<width,integer>

0101.1011101010

integer     fractional

width



particle 0   particle 1   particle 2   . . .   particle 9

(6 features/particle)

(20 features/particle)

(5 features/particle)
(50 features)

(20 features)

(10 features)

(1 feature)
b-tag score

hls 4 ml

Pointwise convolution
(per particle dense layer)

Dense layer

FPGA AUC / Expected AUC

Deregionizer
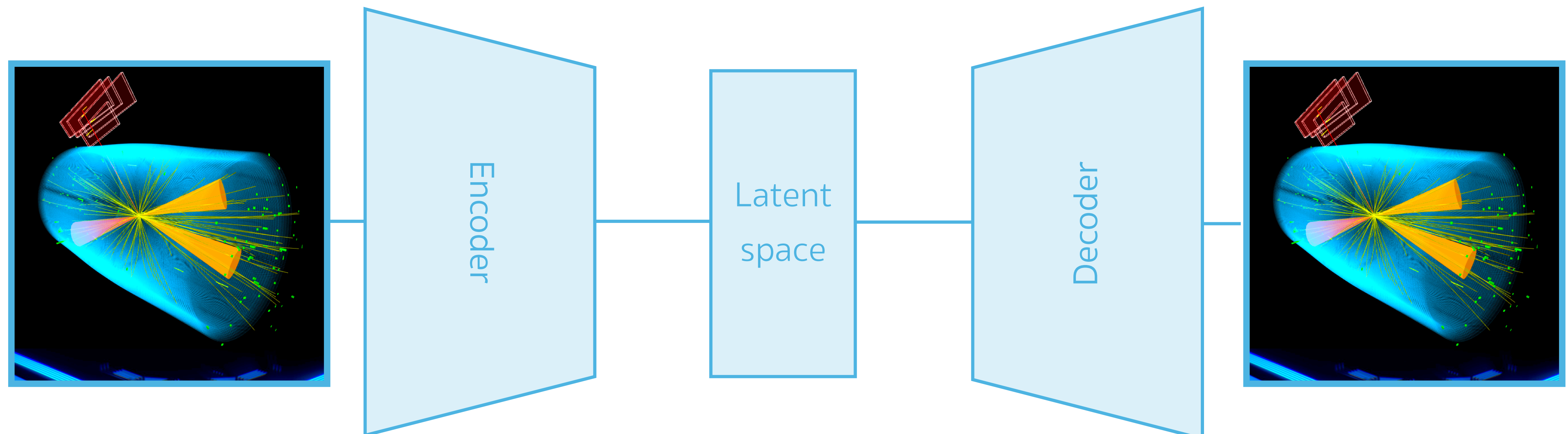
Jet loop

Jet sort
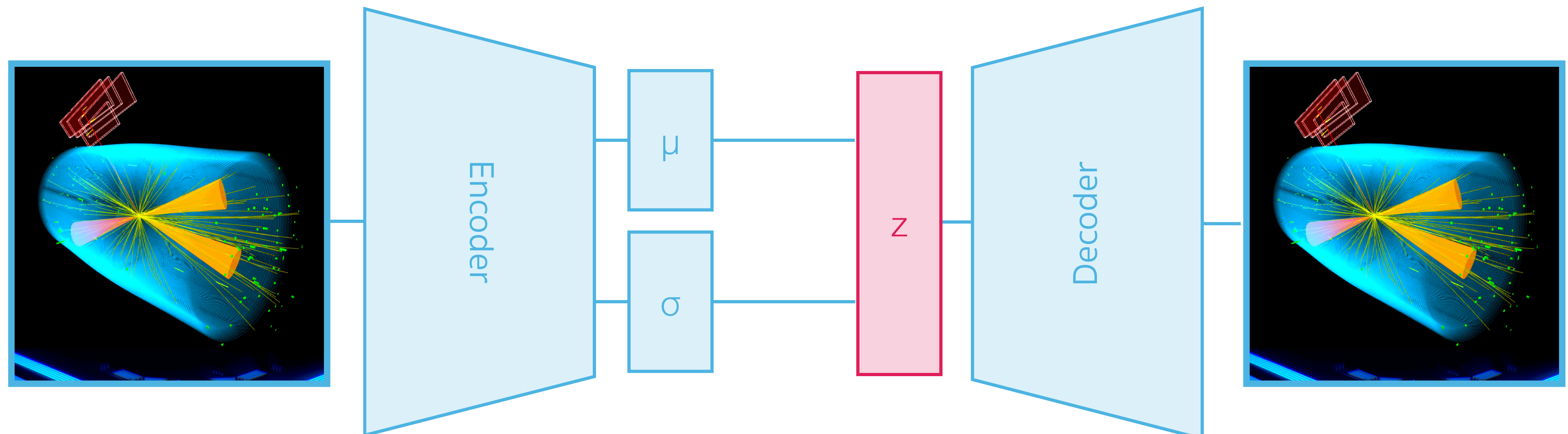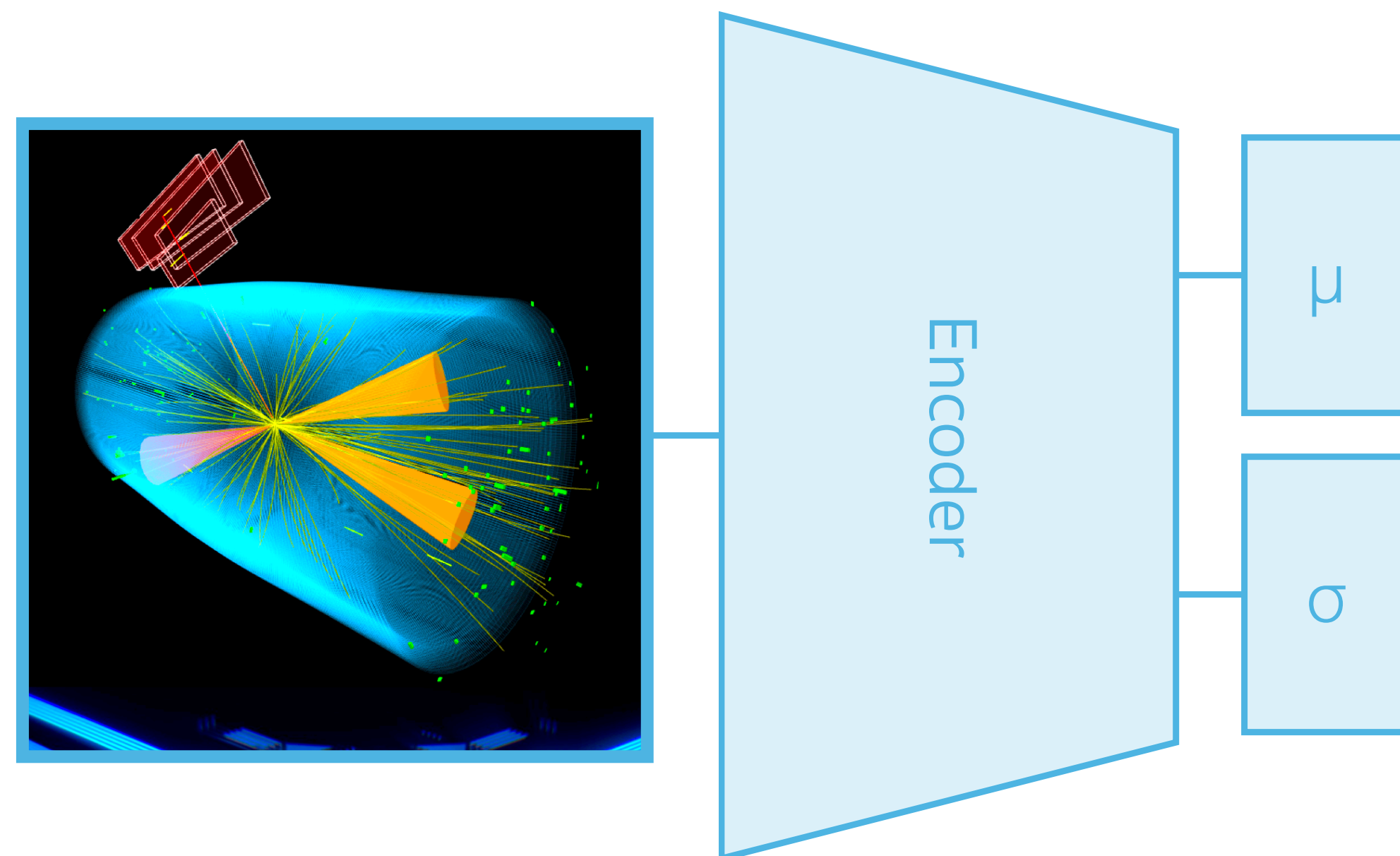
b-tag NN

<25,17> <30,22> <35,27>
fixed-point precision

▸ Challenge: if new physics has an unexpected signature that doesn't align with existing triggers, precious BSM events may be discarded at trigger level

▸ Challenge: if new physics has an unexpected signature that doesn't align with existing triggers, precious BSM events may be discarded at trigger level

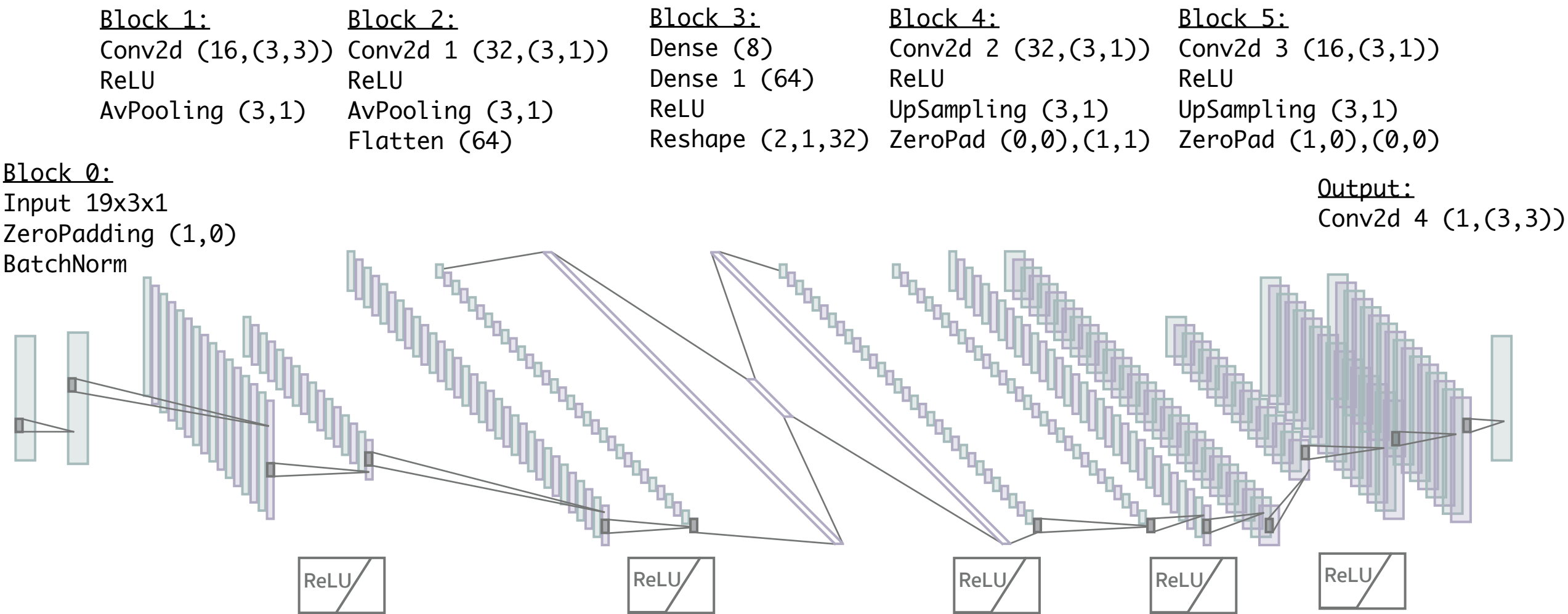▸ Can we use unsupervised algorithms to detect non-SM-like anomalies?

▸ Challenge: if new physics has an unexpected signature that doesn't align with existing triggers, precious BSM events may be discarded at trigger level

▸ Can we use unsupervised algorithms to detect non-SM-like anomalies?

  ▸ Autoencoders (AEs): compress input to a smaller dimensional latent space then decompress and calculate difference



Figure 4: Distribution of four anomaly detection scores (IO AD for AE and VAE models, $R_z$ and $D_{KL}$ ADs for the VAE models) for the DNN model, for the SM cocktail and the four new physics benchmark models.

The model performance is assessed using the four new physics benchmark models. The (ROC) curves in Fig. 6 show the dependence of the true positive rate (TPR) as a function computing by changing the lower threshold applied on the different anomaly scores. We detection performance quoting the area under the ROC curve (AUC) and the TPR corre of SM false positive rate $\varepsilon_{SM} = 10^{-5}$ (see Table 1), which corresponds to the average every month[1].

Encoder

Latent space

Decoder

7

▸ Challenge: if new physics has an unexpected signature that doesn't align with existing triggers, precious BSM events may be discarded at trigger level

▸ Can we use unsupervised algorithms to detect non-SM-like anomalies?

  ▸ Autoencoders (AEs): compress input to a smaller dimensional latent space then decompress and calculate difference

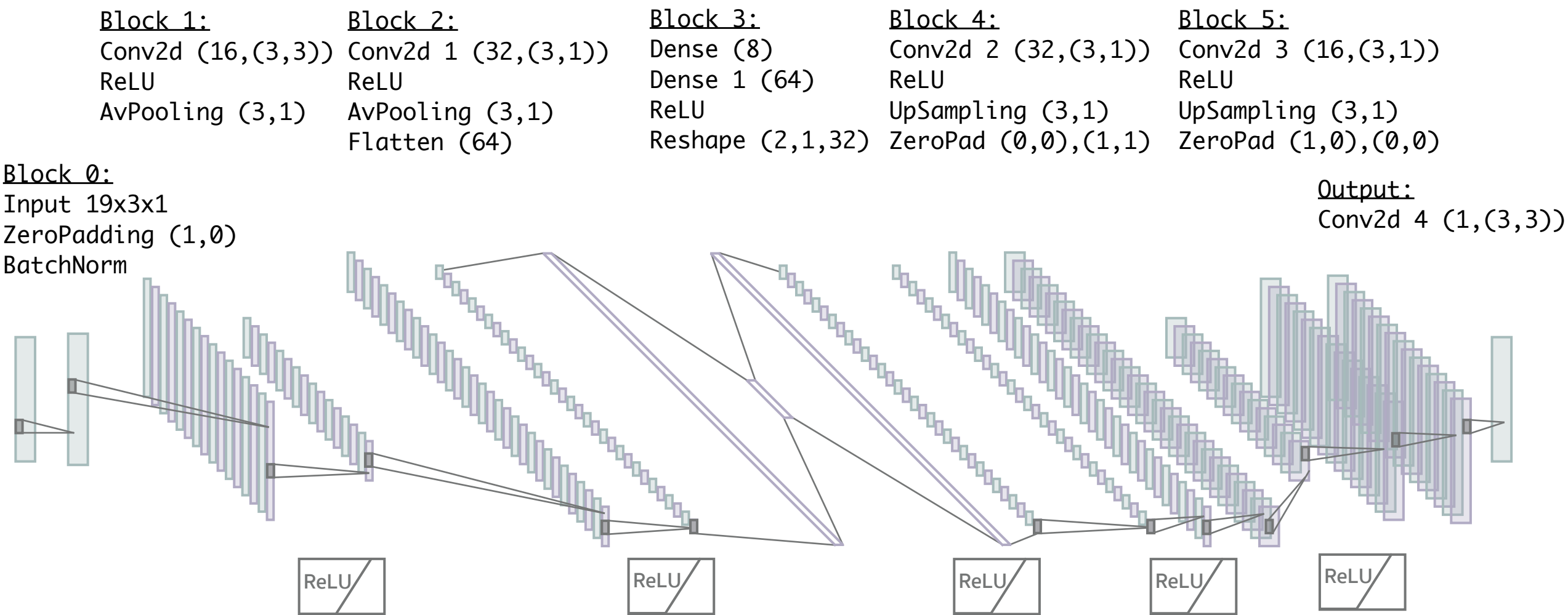  ▸ Variational autoencoders (VAEs): model the latent space as a probability distribution; possible to detect anomalies purely with latent space variables

▸ Challenge: if new physics has an unexpected signature that doesn't align with existing triggers, precious BSM events may be discarded at trigger level

▸ Can we use unsupervised algorithms to detect non-SM-like anomalies?

  ▸ Autoencoders (AEs): compress input to a smaller dimensional latent space then decompress and calculate difference

  ▸ Variational autoencoders (VAEs): model the latent space as a probability distribution; possible to detect anomalies purely with latent space variables

Encoder

μ

σ

Key observation: Can build an anomaly score from the latent space of VAE directly! No need to run decoder!

$$R_z = \sum_i \frac{\mu_i^2}{\sigma_i^2}$$

# FPGA IMPLEMENTATION

Nat. Mach. Intell. 4, 154 (2022)
Data challenge: mpp-hep.github.io/ADC2021    32

Block 1:
Conv2d (16,(3,3))
ReLU
AvPooling (3,1)

Block 2:
Conv2d 1 (32,(3,1))
ReLU
AvPooling (3,1)
Flatten (64)

Block 3:
Dense (8)
Dense 1 (64)
ReLU
Reshape (2,1,32)

Block 4:
Conv2d 2 (32,(3,1))
ReLU
UpSampling (3,1)
ZeroPad (0,0),(1,1)

Block 5:
Conv2d 3 (16,(3,1))
ReLU
UpSampling (3,1)
ZeroPad (1,0),(0,0)

Block 0:
Input 19x3x1
ZeroPadding (1,0)
BatchNorm

Output:
Conv2d 4 (1,(3,3))



ReLU   ReLU   ReLU   ReLU   ReLU

# FPGA IMPLEMENTATION

Nat. Mach. Intell. 4, 154 (2022)
Data challenge: mpp-hep.github.io/ADC2021    32

▸ CNNs as the basis for (V)AEs for anomaly detection

Block 1:
Conv2d (16,(3,3))
ReLU
AvPooling (3,1)

Block 2:
Conv2d 1 (32,(3,1))
ReLU
AvPooling (3,1)
Flatten (64)

Block 3:
Dense (8)
Dense 1 (64)
ReLU
Reshape (2,1,32)

Block 4:
Conv2d 2 (32,(3,1))
ReLU
UpSampling (3,1)
ZeroPad (0,0),(1,1)

Block 5:
Conv2d 3 (16,(3,1))
ReLU
UpSampling (3,1)
ZeroPad (1,0),(0,0)

Block 0:
Input 19x3x1
ZeroPadding (1,0)
BatchNorm

Output:
Conv2d 4 (1,(3,3))

▸ CNNs as the basis for (V)AEs for anomaly detection

▸ Good anomaly detection performance for unseen signals
(LQ → b$\tau$, A → 4l, **h$^\pm$ → $\tau$v**, h$^0$ → $\tau\tau$)



CNN ROC h$^\pm$ → $\tau$v

— IO VAE (AUC = 95%)
— VAE $D_{KL}$ (AUC = 86%)
— VAE $R_z$ (AUC = 86%)
— IO AE (AUC = 96%)

Block 0:
Input 19x3x1
ZeroPadding (1,0)
BatchNorm

Block 1:
Conv2d (16,(3,3))
ReLU
AvPooling (3,1)

Block 2:
Conv2d 1 (32,(3,1))
ReLU
AvPooling (3,1)
Flatten (64)

Block 3:
Dense (8)
Dense 1 (64)
ReLU
Reshape (2,1,32)

Block 4:
Conv2d 2 (32,(3,1))
ReLU
UpSampling (3,1)
ZeroPad (0,0),(1,1)

Block 5:
Conv2d 3 (16,(3,1))
ReLU
UpSampling (3,1)
ZeroPad (1,0),(0,0)

Output:
Conv2d 4 (1,(3,3))

▸ CNNs as the basis for (V)AEs for anomaly detection

▸ Good anomaly detection performance for unseen signals
(LQ → b$\tau$, A → 4l, **h$^\pm$ → $\tau$v**, h$^0$ → $\tau\tau$)

▸ **VAE** fits in latency and resource requirements for HL-LHC!



Block 1:
Conv2d (16,(3,3))
ReLU
AvPooling (3,1)

Block 2:
Conv2d 1 (32,(3,1))
ReLU
AvPooling (3,1)
Flatten (64)

Block 3:
Dense (8)
Dense 1 (64)
ReLU
Reshape (2,1,32)

Block 4:
Conv2d 2 (32,(3,1))
ReLU
UpSampling (3,1)
ZeroPad (0,0),(1,1)

Block 5:
Conv2d 3 (16,(3,1))
ReLU
UpSampling (3,1)
ZeroPad (1,0),(0,0)

Block 0:
Input 19x3x1
ZeroPadding (1,0)
BatchNorm

Output:
Conv2d 4 (1,(3,3))

| Model | DSP [%] | LUT [%] | FF [%] | BRAM [%] | Latency [ns] | II [ns] | AUC [%] | TPR @ FPR=$10^{-5}$ |
|---|---|---|---|---|---|---|---|---|
| CNN VAE R$_z$ | 10 | **12** | **4** | **2** | **365** | **115** | 86 | 0.06% |

CNN ROC h$^\pm$ → $\tau$v
— IO VAE (AUC = 95%)
— VAE D$_{KL}$ (AUC = 86%)
— VAE R$_z$ (AUC = 86%)
— IO AE (AUC = 96%)

▸ CNNs as the basis for (V)AEs for anomaly detection

▸ Good anomaly detection performance for unseen signals (LQ → b$\tau$, A → 4l, **h$^\pm$ → $\tau$v**, h$^0$ → $\tau\tau$)

▸ **VAE** fits in latency and resource requirements for HL-LHC!

*Stay tuned for Run 3…*



CNN ROC h$^\pm$ → $\tau$v

— IO VAE (AUC = 95%)

— VAE $D_{KL}$ (AUC = 86%)

— VAE $R_z$ (AUC = 86%)

— IO AE (AUC = 96%)

Block 0:
Input 19x3x1
ZeroPadding (1,0)
BatchNorm

Block 1:
Conv2d (16,(3,3))
ReLU
AvPooling (3,1)

Block 2:
Conv2d 1 (32,(3,1))
ReLU
AvPooling (3,1)
Flatten (64)

Block 3:
Dense (8)
Dense 1 (64)
ReLU
Reshape (2,1,32)

Block 4:
Conv2d 2 (32,(3,1))
ReLU
UpSampling (3,1)
ZeroPad (0,0),(1,1)

Block 5:
Conv2d 3 (16,(3,1))
ReLU
UpSampling (3,1)
ZeroPad (1,0),(0,0)

Output:
Conv2d 4 (1,(3,3))

| Model | DSP [%] | LUT [%] | FF [%] | BRAM [%] | Latency [ns] | II [ns] | AUC [%] | TPR @ FPR=$10^{-5}$ |
|---|---|---|---|---|---|---|---|---|
| CNN VAE $R_z$ | 10 | **12** | **4** | **2** | **365** | **115** | 86 | 0.06% |

▸ Traditional tracking algorithms scale quadratically with the number of hits

▸ Traditional tracking algorithms scale quadratically with the number of hits

▸ New algorithms (based on **graph neural networks**) may be able to do better

▸ Traditional tracking algorithms scale quadratically with the number of hits

▸ New algorithms (based on **graph neural networks**) may be able to do better

▸ Proof of concept study: use GNN to classify good track segments (edges)

▸ Traditional tracking algorithms scale quadratically with the number of hits

▸ New algorithms (based on **graph neural networks**) may be able to do better

▸ Proof of concept study: use GNN to classify good track segments (edges)

  ▸ Can this fit on an FPGA?

Sector: (1, 0)

Sector: (1, 0)

▶ Build realistic (segmented) graphs for L1 trigger applications

▸ Build realistic (segmented) graphs for L1 trigger applications

▸ ≤8-bit quantized GNN can achieve good edge classification performance

▸ Small graphs (~30 nodes, ~60 edges) easily fit on 1 FPGA

▸ Small graphs (~30 nodes, ~60 edges) easily fit on 1 FPGA

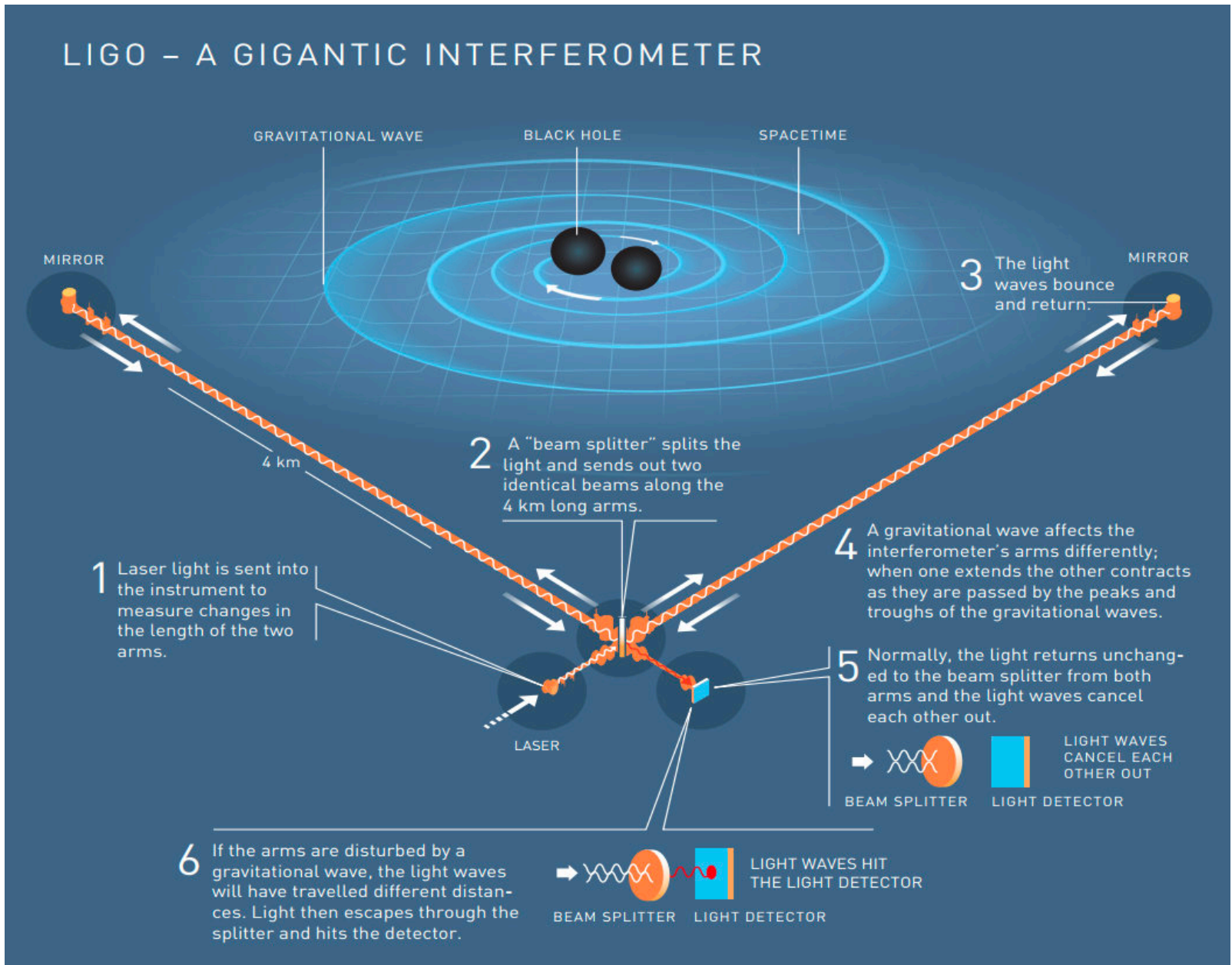▸ Within L1T latency (300 ns) and II (50 ns) requirements

Throughput-optimized, logic synth.
28 nodes, 56 edges
RF = 8

- LUT
- DSP
- FF
- BRAM

Throughput-optimized, C synth.
28 nodes, 56 edges
RF = 8

- Latency
- II

(1 clock cycle = 5 ns)

Similar algorithms for $\tau \to 3\mu$ @ LHC and tracking in sPHENIX

pileup

$\tau$

$\mu$

ICLR 2023

▸ Small graphs (~30 nodes, ~60 edges) easily fit on 1 FPGA

▸ Within L1T latency (300 ns) and II (50 ns) requirements

ECML PKDD 2022

Design to Simulation

Inner Tracker (MVTX and INTT)

MVTX Ladders modeled in details

R ~ 4 cm

sPhenix Trigger Detector along the beam pipe

▶ Nuclear physics

Design to Simulation

Inner Tracker (MVTX and INTT)

MVTX Ladders modeled in details

R ~ 4 cm

sPhenix Trigger Detector along the beam pipe

- Nuclear physics
- Accelerator control



Design to Simulation

Inner Tracker (MVTX and INTT)

MVTX Ladders modeled in details

R ~ 4 cm

sPhenix Trigger Detector along the beam pipe



Booster ν beam
*MicroBooNE, ICARUS, SBND*

Booster
proton energy: 8 GeV

NuMI ν beam
*NOvA experiment*

Main Injector
proton energy: 120 GeV

DUNE ν beam

- Nuclear physics

- Accelerator control

- Neutrino physics



Design to Simulation

Inner Tracker (MVTX and INTT)

MVTX Ladders modeled in details

R ~ 4 cm

sPhenix Trigger Detector along the beam pipe



Booster ν beam
MicroBooNE, ICARUS, SBND

Booster
proton energy: 8 GeV

NuMI ν beam
NOvA experiment

Main Injector
proton energy: 120 GeV

DUNE ν beam

- Nuclear physics
- Accelerator control
- Neutrino physics
- Multi-messenger astronomy



sPhenix Trigger Detector along the beam pipe





LIGO – A GIGANTIC INTERFEROMETER



© 2012 Encyclopædia Britannica, Inc.

- Nuclear physics
- Accelerator control
- Neutrino physics
- Multi-messenger astronomy
- Electron & X-ray microscopy

Design to Simulation

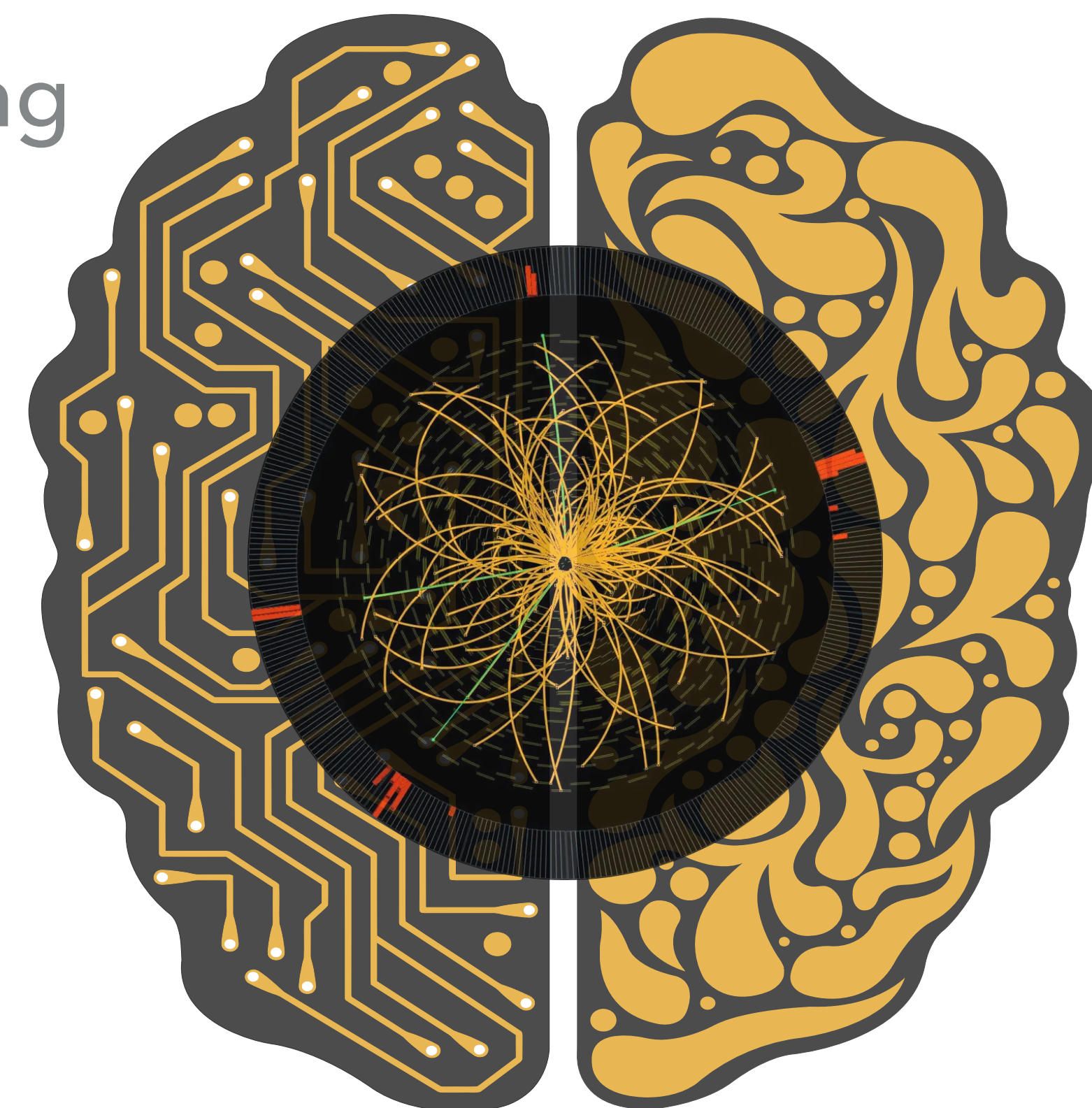Inner Tracker (MVTX and INTT)

MVTX Ladders modeled in details

R ~ 4 cm

sPhenix Trigger Detector along the beam pipe

Booster ν beam
MicroBooNE, ICARUS, SBND

Booster
proton energy: 8 GeV

NuMI ν beam
NOvA experiment

Main Injector
proton energy: 120 GeV

DUNE ν beam

LIGO – A GIGANTIC INTERFEROMETER

GRAVITATIONAL WAVE   BLACK HOLE   SPACETIME

MIRROR   MIRROR

3 The light waves bounce and return.

2 A "beam splitter" splits the light and sends out two identical beams along the 4 km long arms.

4 km

1 Laser light is sent into the instrument to measure changes in the length of the two arms.

4 A gravitational wave affects the interferometer's arms differently; when one extends the other contracts as they are passed by the peaks and troughs of the gravitational waves.

5 Normally, the light returns unchanged to the beam splitter from both arms and the light waves cancel each other out.

LASER

LIGHT WAVES CANCEL EACH OTHER OUT
BEAM SPLITTER   LIGHT DETECTOR

6 If the arms are disturbed by a gravitational wave, the light waves will have travelled different distances. Light then escapes through the splitter and hits the detector.

LIGHT WAVES HIT THE LIGHT DETECTOR
BEAM SPLITTER   LIGHT DETECTOR

electron gun

electron beam

anode

le

magnetic lens

backscattered electron detector

secondary electron detector

specimen

stage

secondary electron detector

stage

© 2012 Encyclopædia Britannica, Inc.

- ▸ Nuclear physics
- ▸ Accelerator control
- ▸ Neutrino physics
- ▸ Multi-messenger astronomy
- ▸ Electron & X-ray microscopy
- ▸ Neuroscience



Design to Simulation

Inner Tracker (MVTX and INTT)

MVTX Ladders modeled in details

R ~ 4 cm

sPhenix Trigger Detector along the beam pipe



Booster ν beam
MicroBooNE, ICARUS, SBND

Booster
proton energy: 8 GeV

NuMI ν beam
NOvA experiment

DUNE ν beam

Main Injector
proton energy: 120 GeV



LIGO – A GIGANTIC INTERFEROMETER

GRAVITATIONAL WAVE      BLACK HOLE      SPACETIME

MIRROR      MIRROR

3  The light waves bounce and return.

2  A "beam splitter" splits the light and sends out two identical beams along the 4 km long arms.

4 km

4  A gravitational wave affects the interferometer's arms differently; when one extends the other contracts as they are passed by the peaks and troughs of the gravitational waves.

1  Laser light is sent into the instrument to measure changes in the length of the two arms.

5  Normally, the light returns unchanged to the beam splitter from both arms and the light waves cancel each other out.

LASER

LIGHT WAVES CANCEL EACH OTHER OUT
BEAM SPLITTER      LIGHT DETECTOR

6  If the arms are disturbed by a gravitational wave, the light waves will have travelled different distances. Light then escapes through the splitter and hits the detector.

LIGHT WAVES HIT THE LIGHT DETECTOR
BEAM SPLITTER      LIGHT DETECTOR



electron gun

electron beam

anode

magnetic lens

backscattered electron detector

secondary electron detector

specimen

stage

1 mm

▸ ML allows us to better reconstruct our data and save potentially overlooked data

▸ ML allows us to better reconstruct our data and save potentially overlooked data

▸ **Codesign** principles can enable ML on hardware with stringent constraints

▸ ML allows us to better reconstruct our data and save potentially overlooked data

▸ **Codesign** principles can enable ML on hardware with stringent constraints

▸ Community (fastmachinelearning.org, e-group hls-fml@cern.ch) and Institute (a3d3.ai) developing open-source tools and techniques to enable this

   ▸ hls4ml: expanding open-source toolkit for translating ML into hardware aimed at trigger applications and more…

▸ ML allows us to better reconstruct our data and save potentially overlooked data

▸ **Codesign** principles can enable ML on hardware with stringent constraints

▸ Community (fastmachinelearning.org, e-group hls-fml@cern.ch) and Institute (a3d3.ai) developing open-source tools and techniques to enable this

  ▸ hls4ml: expanding open-source toolkit for translating ML into hardware aimed at trigger applications and more…

▸ Applications range from momentum regression, to b-tagging, tracking, and more!

  ▸ Enhance **future particle physics program**

Accelerated AI
Algorithms for
Data-Driven
Discovery

BACKUP

Unstructured Pruning

▸ Unstructured pruning: removing some connections regardless of placement

Unstructured Pruning

▸ Unstructured pruning: removing some connections regardless of placement

▸ Structured pruning: removing all input/output connections of particular nodes



Unstructured Pruning

Structured Pruning

▶ Excellent overview talks for reference

▸ **Tools**

    ▸ Accessible workflows like HLS to make hardware more accessible domain scientists

▸ **ML techniques**

    ▸ Efficient training and implementation methods codesigned for specific hardware

▸ **Hardware**

    ▸ Evolving compute platforms, e.g. power-law growth in FPGA logic

**Scientific domain applications**

**Science data pipelines**

**Compute elements & systems**

**Domain-inspired ML**

**ML-specific systems**

**Machine learning techniques**

▸ Reconstruct all events and reject 98% of them in ~12.5 μs

  ▸ Individual algorithms usually have to be < 1 μs and keep up with new events every 25 ns

▸ Latency necessitates all **FPGA** design (many algorithms running on 729 FPGAs!)

  ▸ Individual algorithms usually have to fit on < 1 FPGA

▸ Reconstruct all events and reject 98% of them in ~12.5 µs

  ▸ Individual algorithms usually have to be < 1 µs and keep up with new events every 25 ns

▸ Latency necessitates all **FPGA** design (many algorithms running on 729 FPGAs!)

  ▸ Individual algorithms usually have to fit on < 1 FPGA

▸ Reconstruct all events and reject 98% of them in ~12.5 μs

  ▸ Individual algorithms usually have to be < 1 μs and keep up with new events every 25 ns

▸ Latency necessitates all **FPGA** design (many algorithms running on 729 FPGAs!)

  ▸ Individual algorithms usually have to fit on < 1 FPGA

$$\ell_j^k = \phi \left( \sum_i w_{ij} \ell_i^{k-1} + b_j \right)$$

$$\ell_j^k = \phi \left( \sum_i w_{ij} \ell_i^{k-1} + b_j \right)$$

multiplication

activation function

addition



$\ell_i^{k-1}$

$w_{ij}$

$\ell_j^k$

multiplication

$$\ell_j^k = \phi \left( \sum_i w_{ij} \ell_i^{k-1} + b_j \right)$$

activation function

addition

$\ell_i^{k-1}$

$w_{ij}$

$\ell_j^k$

Maps nicely onto FPGA resources: high I/O, DSPs, LUTs, etc.

▸ Operations can be implemented with core operations (gates)



AND OR XOR

NAND NOR XNOR

▸ Operations can be implemented with core operations (gates)

### LUT

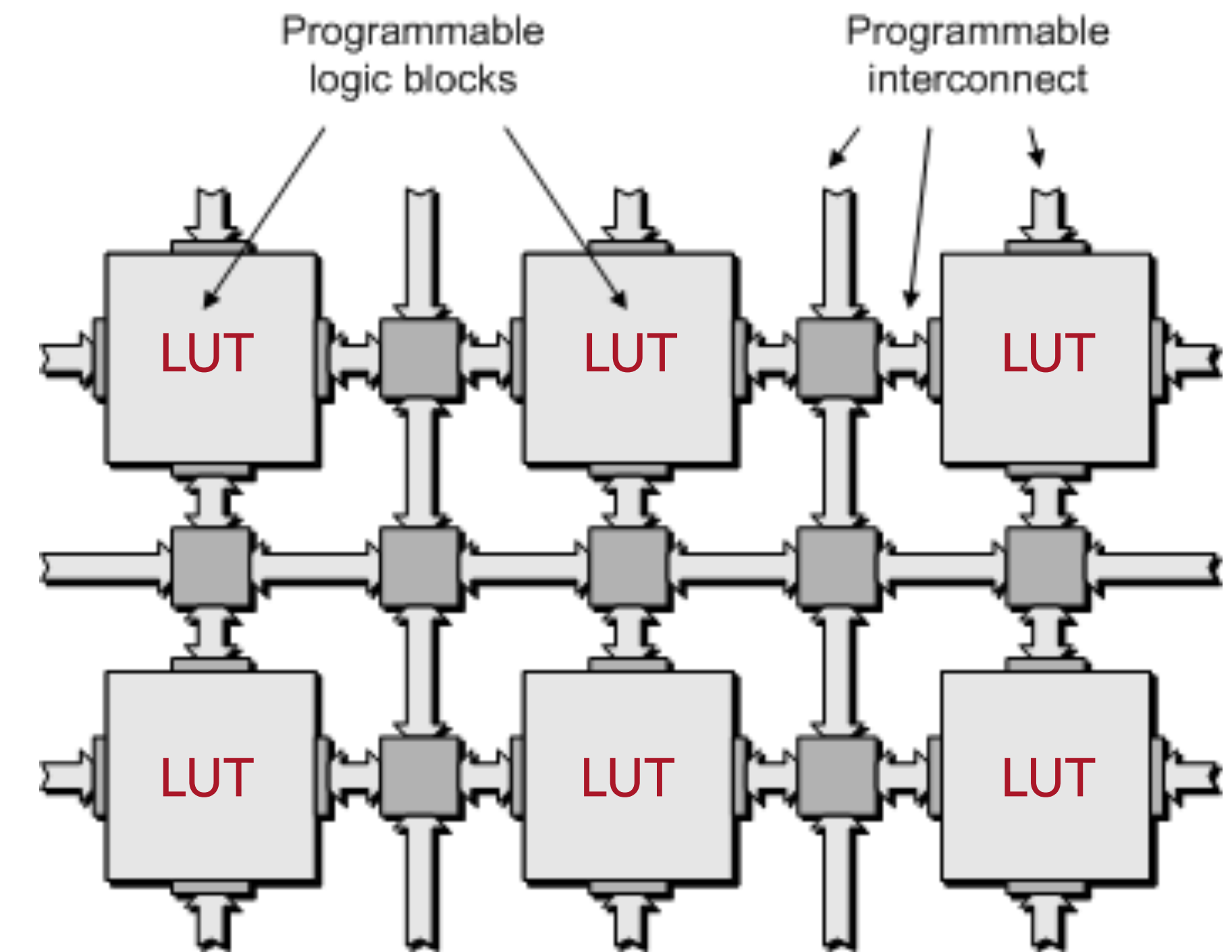| A | B | Output |
|---|---|--------|
| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 1 |

### LUT

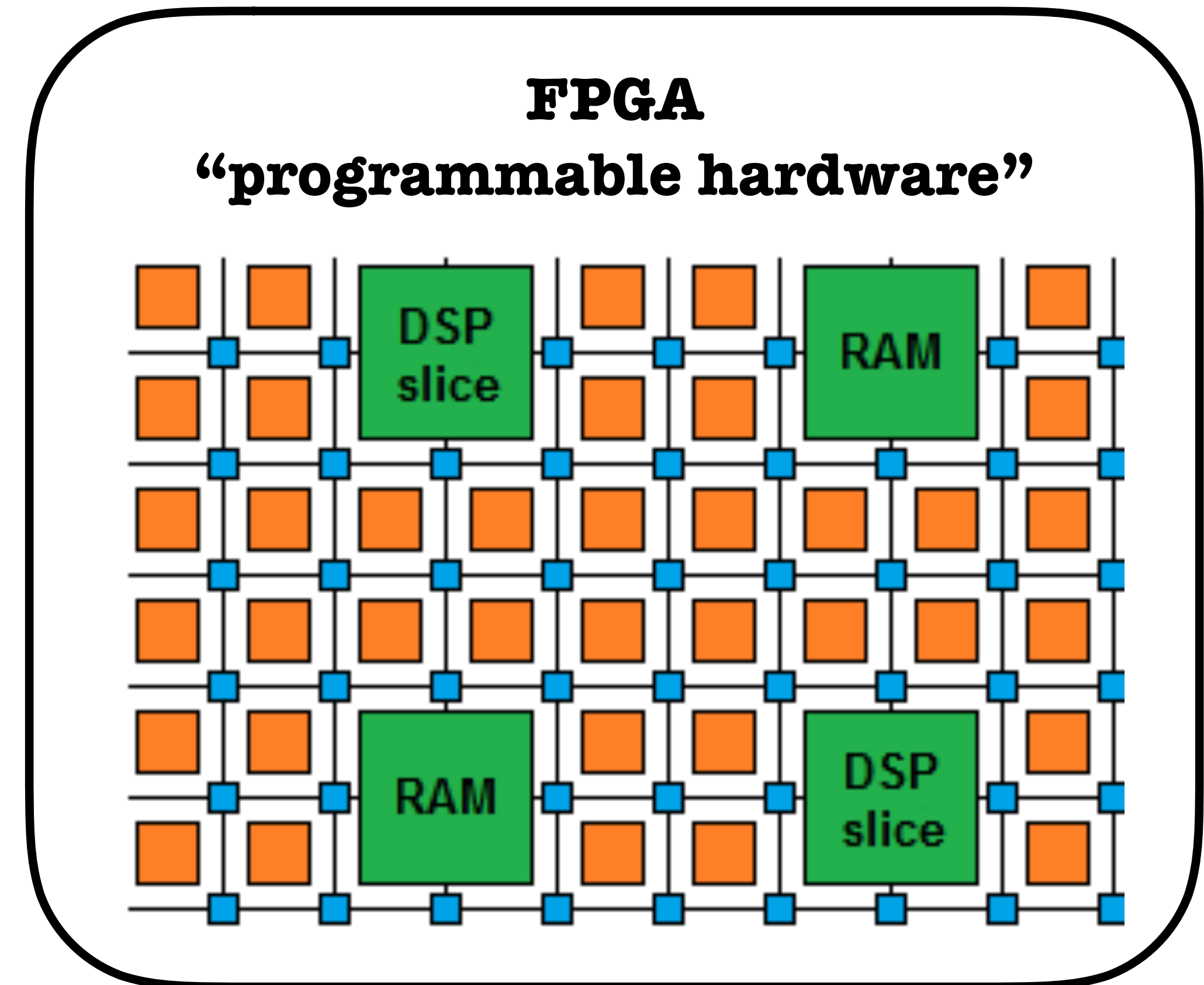| A | B | Output |
|---|---|--------|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |

### LUT

| A | B | Output |
|---|---|--------|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

### LUT

| A | B | Output |
|---|---|--------|
| 0 | 0 | 1 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

### LUT

| A | B | Output |
|---|---|--------|
| 0 | 0 | 1 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 0 |

### LUT

| A | B | Output |
|---|---|--------|
| 0 | 0 | 1 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 1 |

▸ Gates are like look-up tables (LUTs)

▸ Operations can be implemented with core operations (gates)

LUT

| A | B | Output |
|---|---|--------|
| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 1 |

LUT

| A | B | Output |
|---|---|--------|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |

LUT

| A | B | Output |
|---|---|--------|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

LUT

| A | B | Output |
|---|---|--------|
| 0 | 0 | 1 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

LUT

| A | B | Output |
|---|---|--------|
| 0 | 0 | 1 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 0 |

LUT

| A | B | Output |
|---|---|--------|
| 0 | 0 | 1 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 1 |



▸ Gates are like look-up tables (LUTs)

▸ If we can (re-)program arbitrary LUTs and (re-)connect them however we want, we can (re-)implement whatever algorithm we want!

▸ Pros:

    ▸ Reprogrammable interconnects between embedded components that perform multiplication (DSPs), apply logical functions (LUTs), or store memory (BRAM)

    ▸ High throughput I/O: O(100) optical transceivers running at O(15) Gbps

    ▸ Massively parallel

    ▸ Low power

▸ Cons:

    ▸ Requires domain knowledge to program (using VHDL/Verilog)
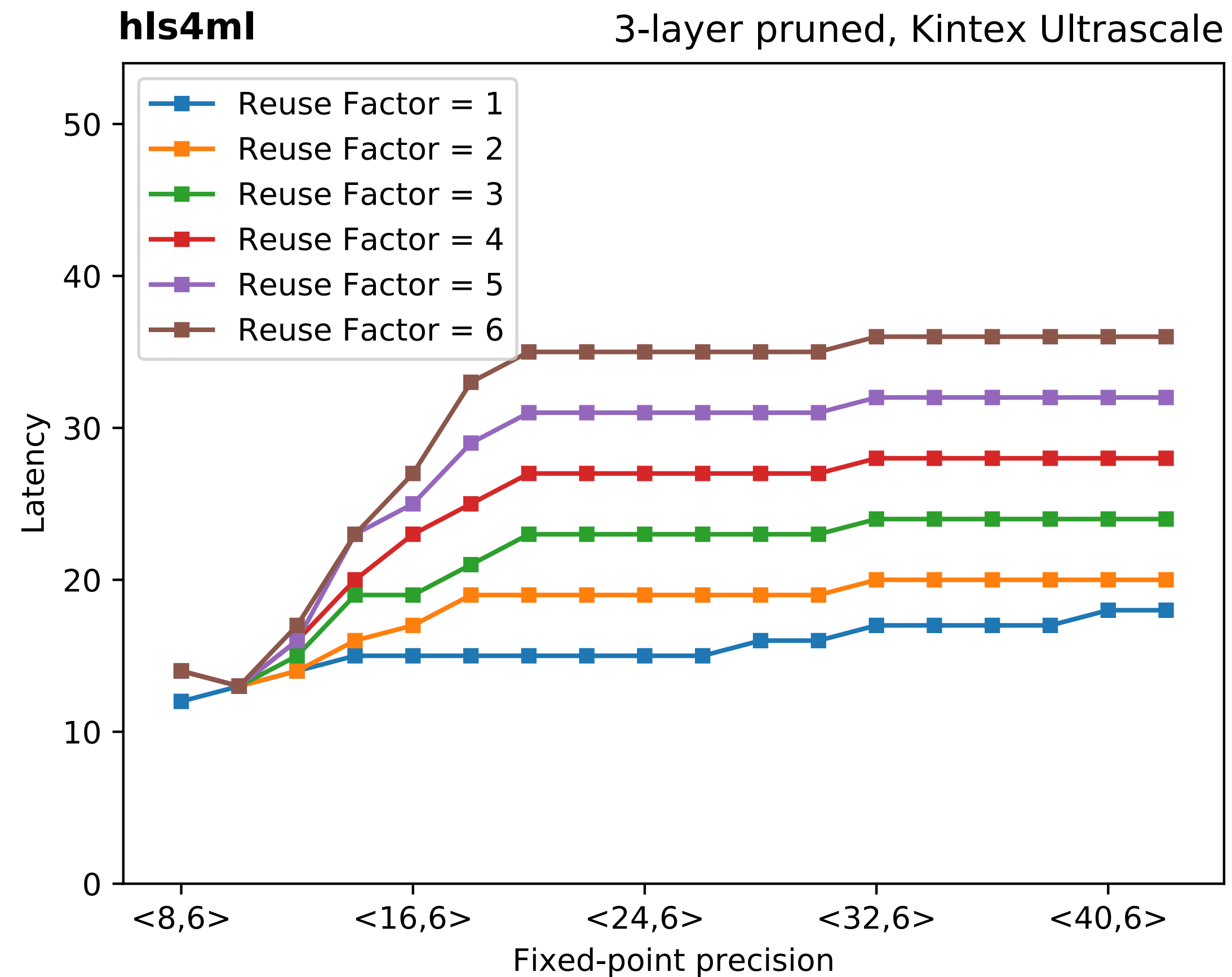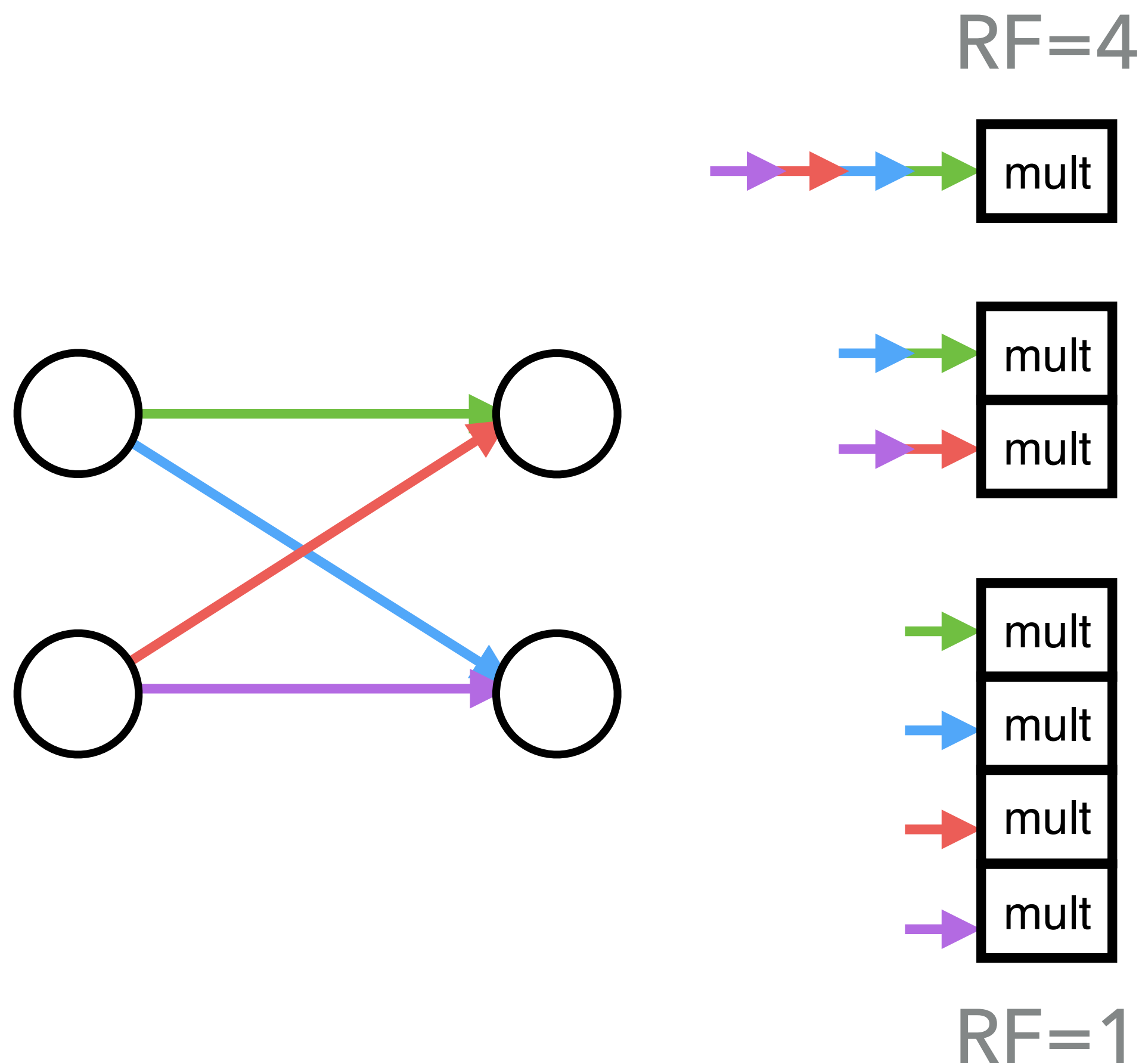
**FPGA**
**"programmable hardware"**

▸ Decreasing reuse factor, increases parallelization and decreases latency



RF=4

RF=1

hls4ml

3-layer pruned, Kintex Ultrascale

Reuse Factor = 1
Reuse Factor = 2
Reuse Factor = 3
Reuse Factor = 4
Reuse Factor = 5
Reuse Factor = 6

Latency

Fixed-point precision

~35 clocks @ 200 MHz = 175 ns

~15 clocks @ 200 MHz = 75 ns

▸ Decreasing reuse factor, increases parallelization and decreases latency



RF=4

RF=1

hls4ml    3-layer pruned, Kintex Ultrascale

~35 clocks @ 200 MHz = 175 ns

~15 clocks @ 200 MHz = 75 ns

▸ Algorithm comfortably fits in latency requirements (<1 μs)

**Fast inference of deep neural networks in FPGAs for particle physics**

J. Duarte,[a] S. Han,[b] P. Harris,[b] S. Jindariani,[a] E. Kreinar,[c] B. Kreis,[a] J. Ngadiuba,[d] M. Pierini,[d] R. Rivera,[a] N. Tran[a,1] and Z. Wu[e]

## Autoencoders on field-programmable gate arrays for real-time, unsupervised new physics detection at 40 MHz at the Large Hadron Collider

Ekaterina Govorkova [1✉], Ema Puljak [1], Thea Aarrestad [1], Thomas James[1], Vladimir Loncar [1,2], Maurizio Pierini [1], Adrian Alan Pol [1], Nicolò Ghielmetti[1,3], Maksymilian Graczyk[1,4], Sioni Summers[1], Jennifer Ngadiuba [5,6], Thong Q. Nguyen[6], Javier Duarte[7] and Zhenbin Wu[8]

MACHINE LEARNING
Science and Technology

### Compressing deep neural networks on FPGAs to binary and ternary precision with `hls4ml`

Jennifer Ngadiuba[1] [©], Vladimir Loncar[1], Maurizio Pierini[1], Sioni Summers[1], Giuseppe Di Guglielmo[2], Javier Duarte[3] [©], Philip Harris[4], Dylan Rankin[4], Sergo Jindariani[3], Mia Liu[5], Kevin Pedro[5], Nhan Tran[3], Edward Kreinar[6], Sheila Sagear[7], Zhenbin Wu[8] and Duc Hoang[9]

## A Reconfigurable Neural Network ASIC for Detector Front-End Data Compression at the HL-LHC

Giuseppe Di Guglielmo [©], Farah Fahim [©], *Member, IEEE*, Christian Herwig [©], Manuel Blanco Valentin, Javier Duarte [©], Cristian Gingu, *Member, IEEE*, Philip Harris, James Hirschauer [©], Martin Kwok, Vladimir Loncar, Yingyi Luo, Llovizna Miranda, Jennifer Ngadiuba, Daniel Noonan, Seda Ogrenci-Memik, Maurizio Pierini, Sioni Summers, and Nhan Tran [©]

## Distance-Weighted Graph Neural Networks on FPGAs for Real-Time Particle Reconstruction in High Energy Physics

### hls4ml: An Open-Source Codesign Workflow to Empower Scientific Low-Power Machine Learning Devices

Farah Fahim*
Benjamin Hawks
Christian Herwig
James Hirschauer
Sergo Jindariani
Nhan Tran*
Fermilab
Batavia, IL, USA

Luca P. Carloni
Giuseppe Di Guglielmo
Columbia University
New York, NY, USA

Philip Harris
Jeffrey Krupa
Dylan Rankin
MIT
Cambridge, MA, USA

Manuel Blanco Valentin
Josiah Hester
Yingyi Luo
John Mamish
Seda Orgrenci-Memik
Northwestern University
Evanston, IL, USA

Thea Aarrestad
Hamza Javed
Vladimir Loncar
Maurizio Pierini
Adrian Alan Pol
Sioni Summers
European Organization for Nuclear
Research (CERN)
Geneva, Switzerland

Javier Duarte
UC San Diego
La Jolla, CA, USA
jduarte@ucsd.edu

Scott Hauck
Shih-Chieh Hsu
University of Washington
Seattle, WA, USA

Jennifer Ngadiuba
Caltech
Pasadena, CA, USA

Mia Liu
Purdue University
West Lafayette, IN, USA

Duc Hoang
Rhodes College
Memphis, TN, USA

Edward Kreinar
HawkEye360
Herndon, VA, USA

Zhenbin Wu
University of Illinois at Chicago
Chicago, IL, USA

arXiv:2103.05579v3 [cs.LG] 23 Mar 2021

## ESP4ML: Platform-Based Design of Systems-on-Chip for Embedded Machine Learning

Davide Giri, Kuan-Lin Chiu, Giuseppe Di Guglielmo, Paolo Mantovani and Luca P. Carloni
*Department of Computer Science · Columbia University*, New York
[davide_giri, chiu, giuseppe, paolo, luca]@cs.columbia.edu

**Fast inference of Boosted Decision Trees in FPGAs for particle physics**

S. Summers,[a,1] G. Di Guglielmo,[b] J. Duarte,[c] P. Harris,[d] D. Hoang,[e] S. Jindariani,[f] E. Kreinar,[g] V. Loncar,[a,h] J. Ngadiuba,[a] M. Pierini,[a] D. Rankin,[d] N. Tran[f] and Z. Wu[i]

## Automatic heterogeneous quantization of deep neural networks for low-latency inference on the edge for particle detectors

Claudionor N. Coelho Jr[1], Aki Kuusela[2], Shan Li[2], Hao Zhuang[2], Jennifer Ngadiuba [3], Thea Klaeboe Aarrestad [4✉], Vladimir Loncar[4,5], Maurizio Pierini[4], Adrian Alan Pol [4] and Sioni Summers[4]

## The Phase-2 Upgrade of the CMS Level-1 Trigger
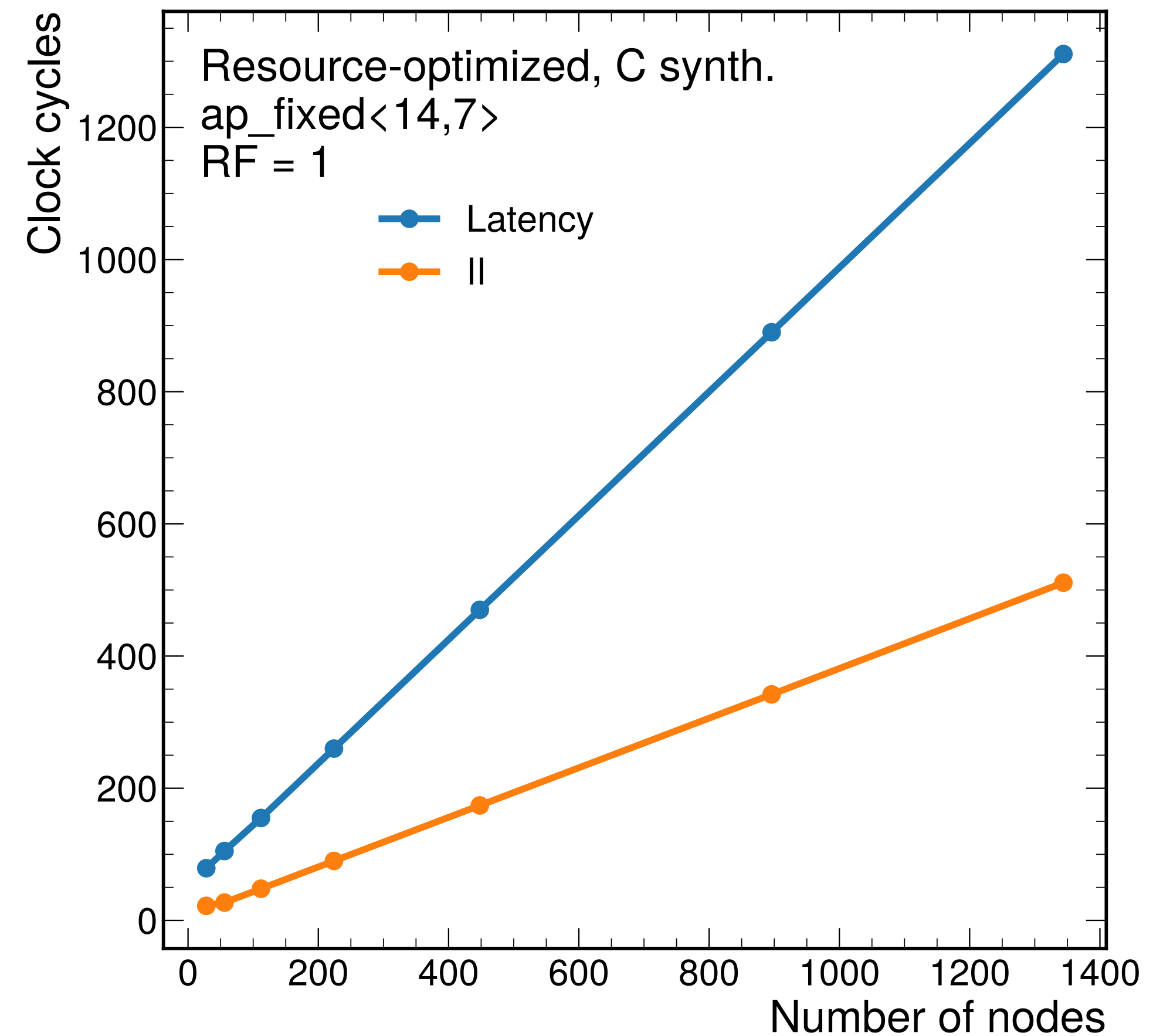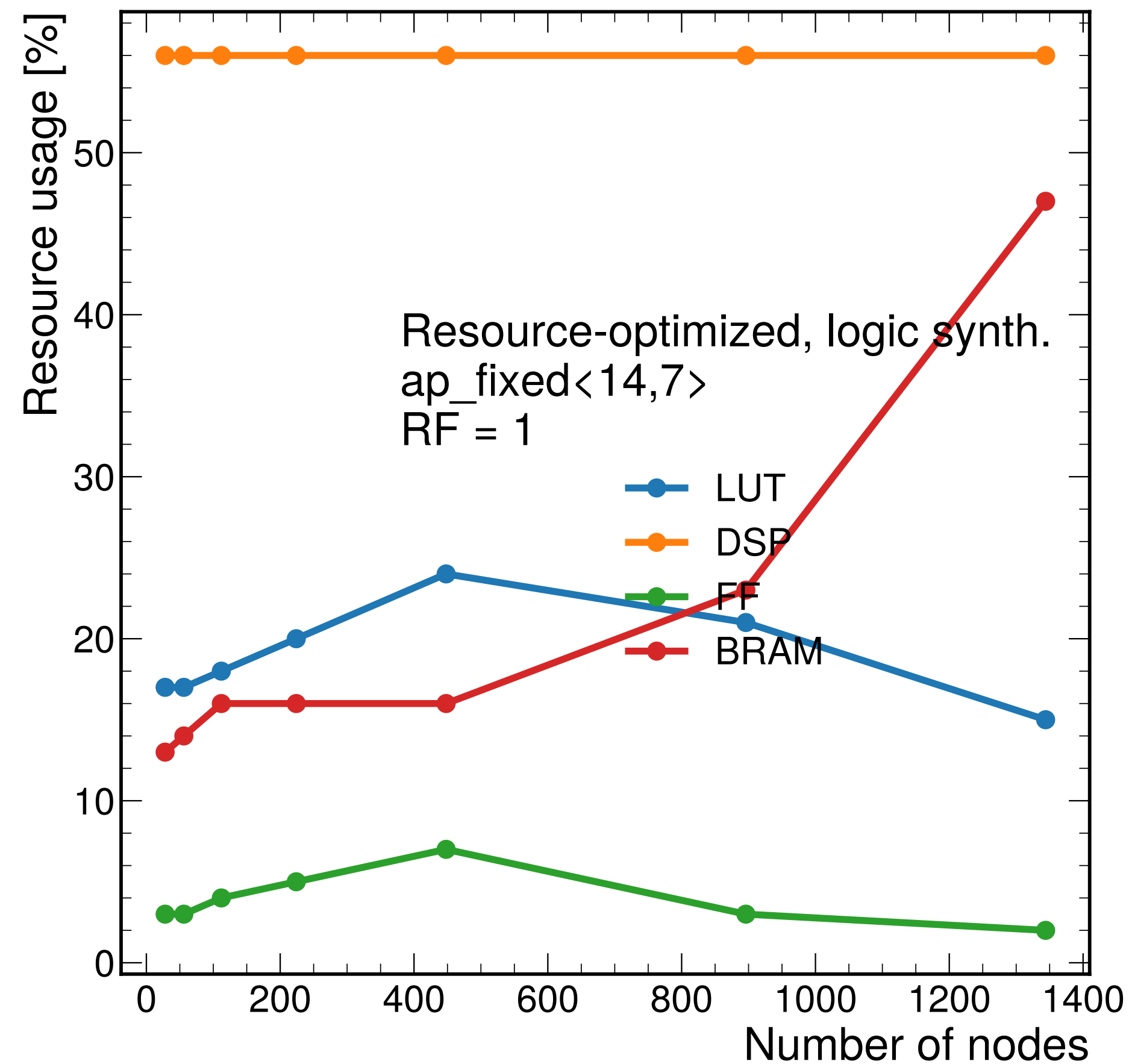### Technical Design Report

## Graph Neural Networks for Charged Particle Tracking on FPGAs

Abdelrahman Elabd[1], Vesal Razavimaleki[2], Shi-Yu Huang[3], Javier Duarte[2*], Markus Atkinson[4], Gage DeZoort[5], Peter Elmer[5], Scott Hauck[6], Jin-Xuan Hu[3], Shih-Chieh Hsu[6,7], Bo-Cheng Lai[3], Mark Neubauer[4*], Isobel Ojalvo[5], Savannah Thais[5] and Matthew Trahms[6]
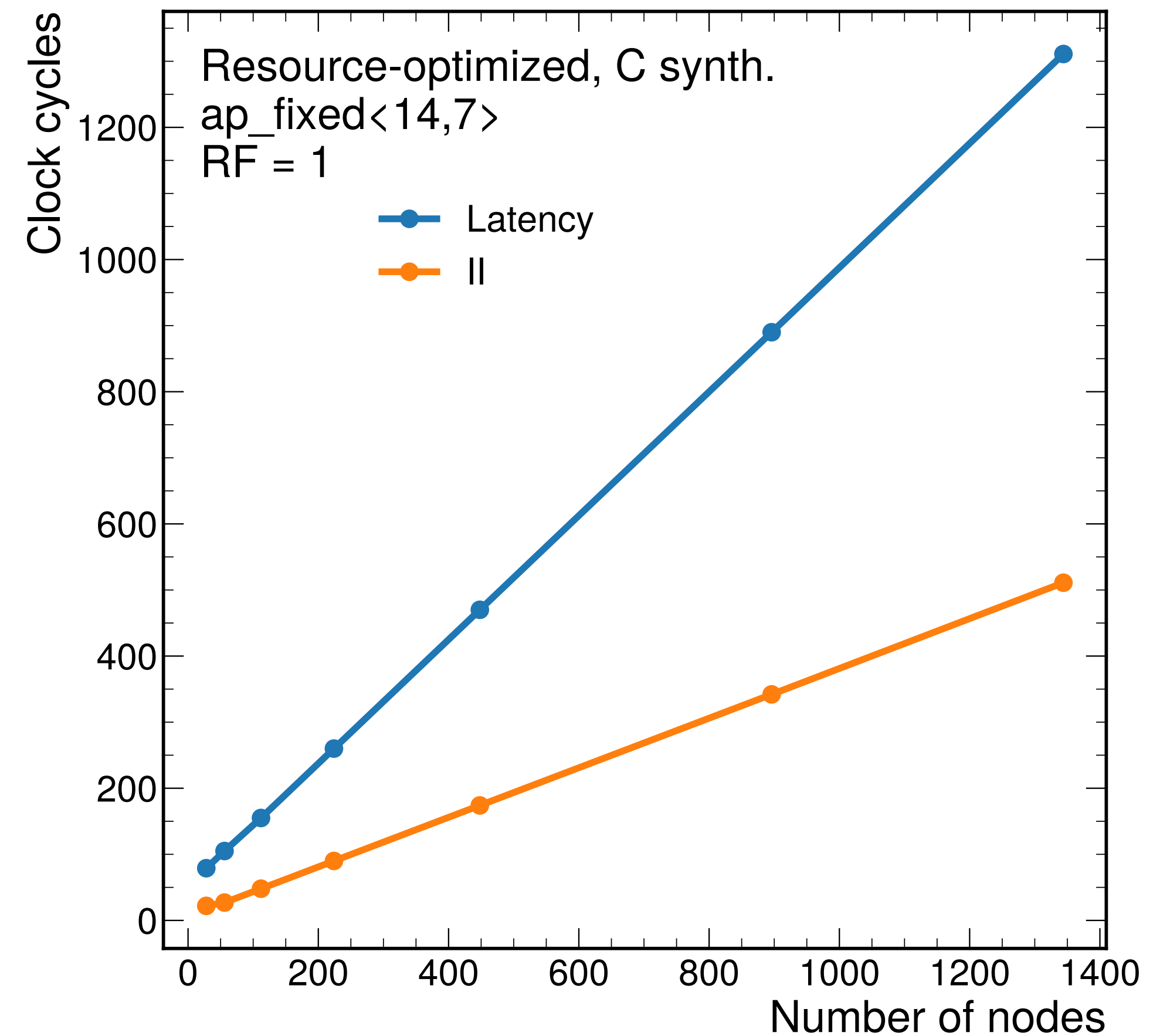
1 clock cycle = 5 ns

1 clock cycle = 5 ns



▸ Modified design can scale to much larger graphs (~1400 nodes, ~2800 edges), for longer latency (6 μs) and II (2 μs)

1. Define generic ML benchmarks for bespoke domain problems that attract interest from a broad community of system and ML experts

2. Design benchmarks to satisfy challenging scientific requirements that overlap with a number of systems

1. Define generic ML benchmarks for bespoke domain problems that attract interest from a broad community of system and ML experts

2. Design benchmarks to satisfy challenging scientific requirements that overlap with a number of systems

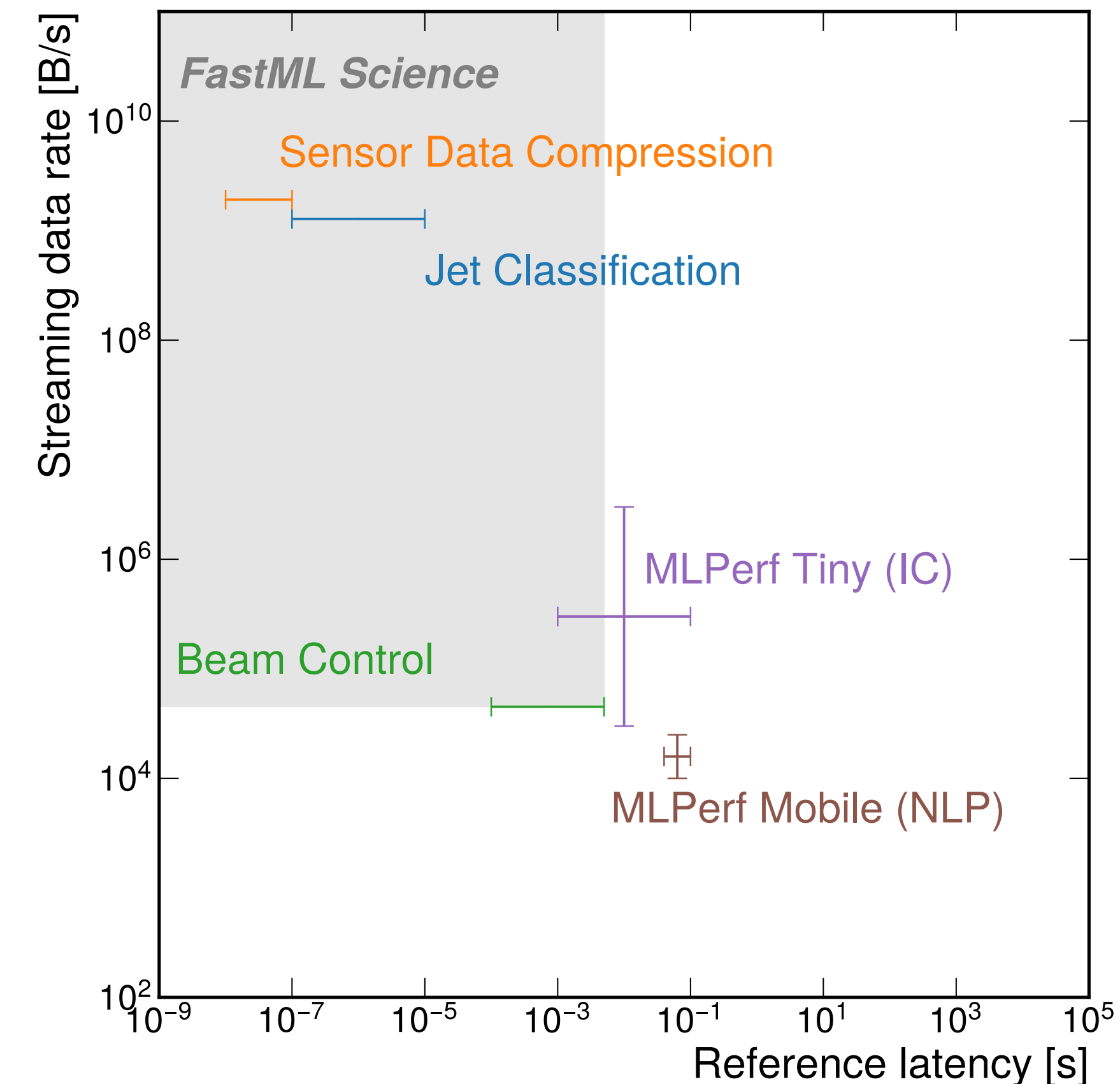### FastML Science Benchmarks: Accelerating Real-Time Scientific Edge Machine Learning

**Javier Duarte** [*][1]  **Nhan Tran** [*][2]  **Ben Hawks** [2]  **Christian Herwig** [2]
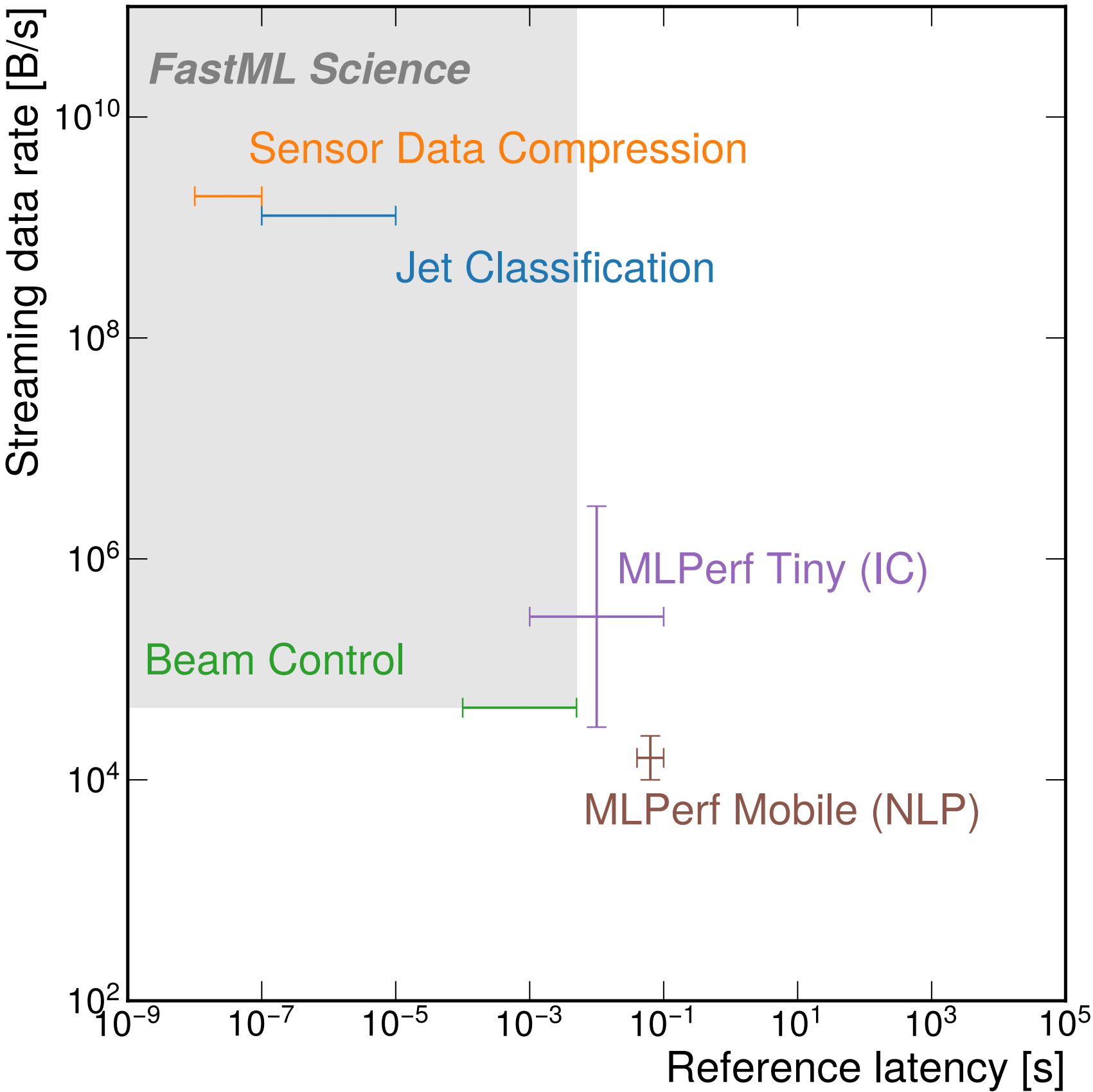**Jules Muhizi** [3]  **Shvetank Prakash** [3]  **Vijay Janapa Reddi** [3]

1. Define generic ML benchmarks for bespoke domain problems that attract interest from a broad community of system and ML experts

2. Design benchmarks to satisfy challenging scientific requirements that overlap with a number of systems

‣ Set of 3 benchmarks inspired by low-latency edge ML use cases in science

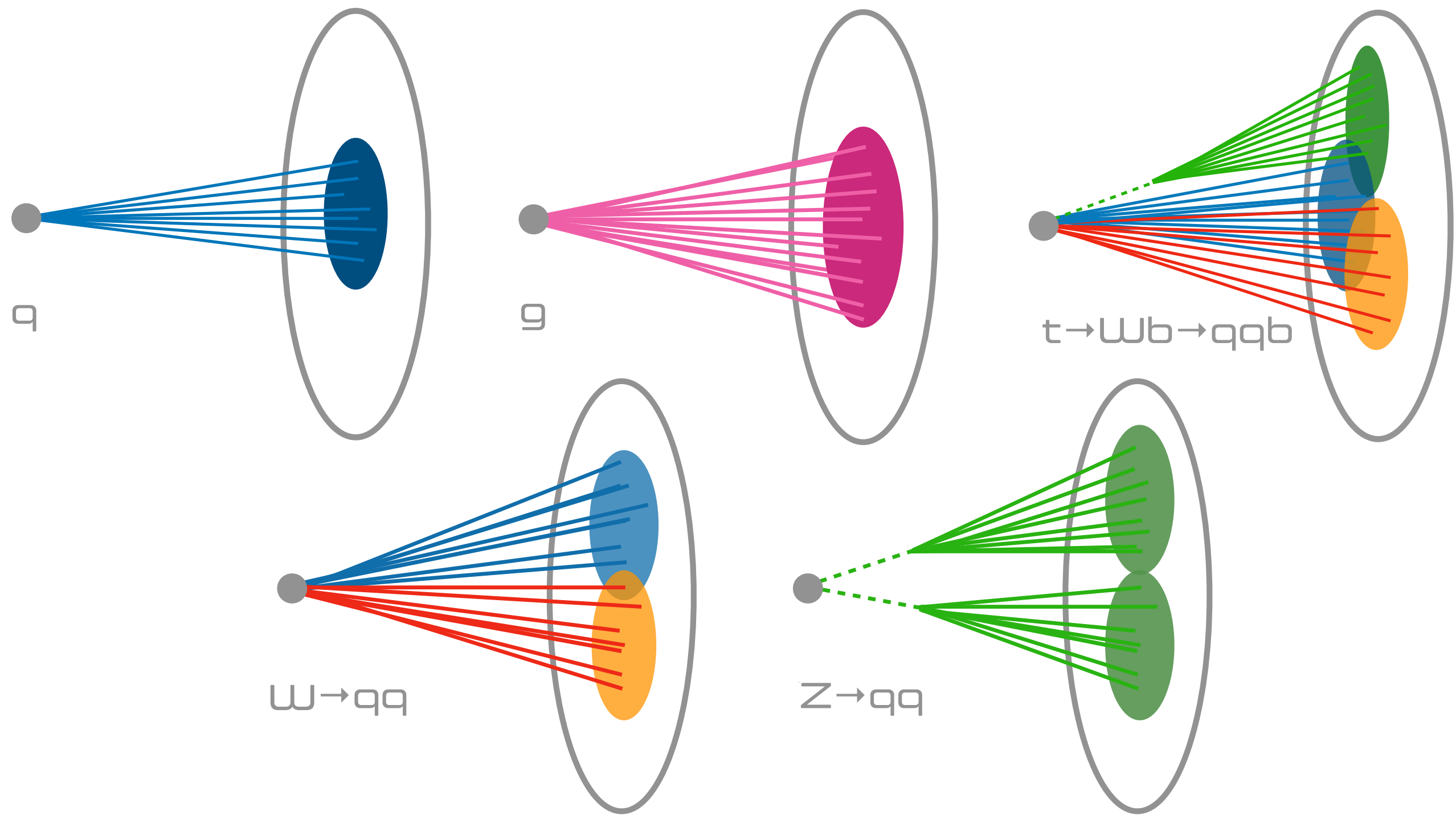‣ Cover a wide range of latency/data rate constraints

**FASTML SCIENCE BENCHMARKS:
ACCELERATING REAL-TIME SCIENTIFIC EDGE MACHINE LEARNING**

**Javier Duarte** [*1]   **Nhan Tran** [*2]   **Ben Hawks** [2]   **Christian Herwig** [2]
**Jules Muhizi** [3]   **Shvetank Prakash** [3]   **Vijay Janapa Reddi** [3]

1. Define generic ML benchmarks for bespoke domain problems that attract interest from a broad community of system and ML experts

2. Design benchmarks to satisfy challenging scientific requirements that overlap with a number of systems

**FᴀsᴛML Sᴄɪᴇɴᴄᴇ Bᴇɴᴄʜᴍᴀʀᴋs:
Aᴄᴄᴇʟᴇʀᴀᴛɪɴɢ Rᴇᴀʟ-Tɪᴍᴇ Sᴄɪᴇɴᴛɪғɪᴄ Eᴅɢᴇ Mᴀᴄʜɪɴᴇ Lᴇᴀʀɴɪɴɢ**

**Javier Duarte** [*,1]  **Nhan Tran** [*,2]  **Ben Hawks** [2]  **Christian Herwig** [2]
**Jules Muhizi** [3]  **Shvetank Prakash** [3]  **Vijay Janapa Reddi** [3]

▸ Set of 3 benchmarks inspired by low-latency edge ML use cases in science

▸ Cover a wide range of latency/data rate constraints

▸ Unique set of qualities

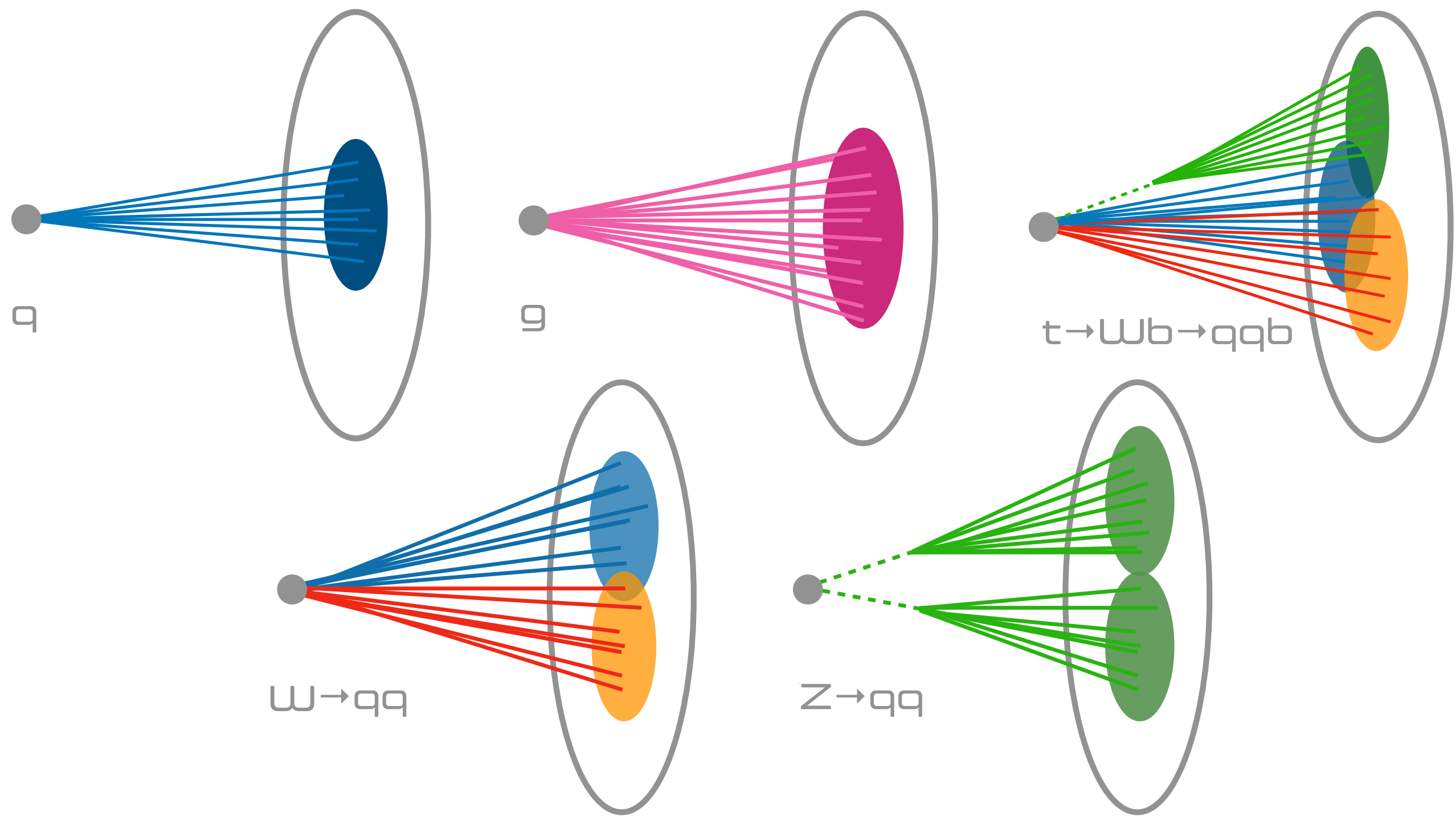| | Formalized Benchmark | Scientific Workload(s) | Edge Computing | Real-Time Constraints |
|---|---|---|---|---|
| **FastML Science Benchmarks (this work)** | ✓ | ✓ | ✓ | ✓ |
| SciMLBench (Thiyagalingam et al., 2021) | ✓ | ✓ | ✓ | x |
| LHC New Physics Dataset (Govorkova et al., 2021) | x | ✓ | ✓ | ✓ |
| MLPerf HPC (Farrell et al., 2021) | ✓ | ✓ | x | x |
| BenchCounil AIBench HPC (BenchCouncil, 2018) | ✓ | ✓ | x | x |
| MLCommons Science (MLCommons, 2020) | ✓ | ✓ | x | x |
| ITU Modulation Classification (ITU, 2021) | x | x | ✓ | ✓ |

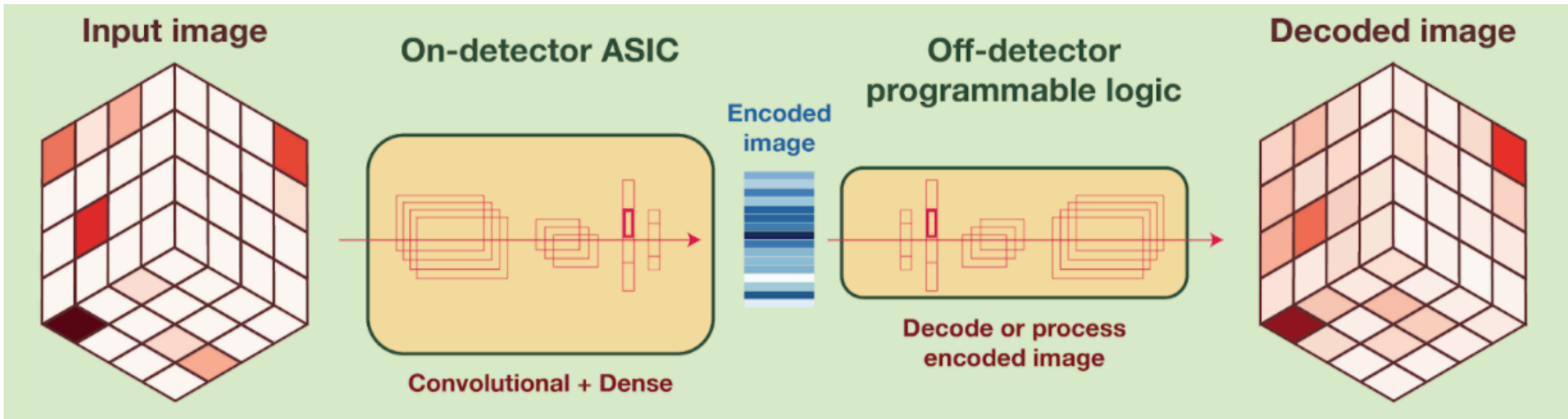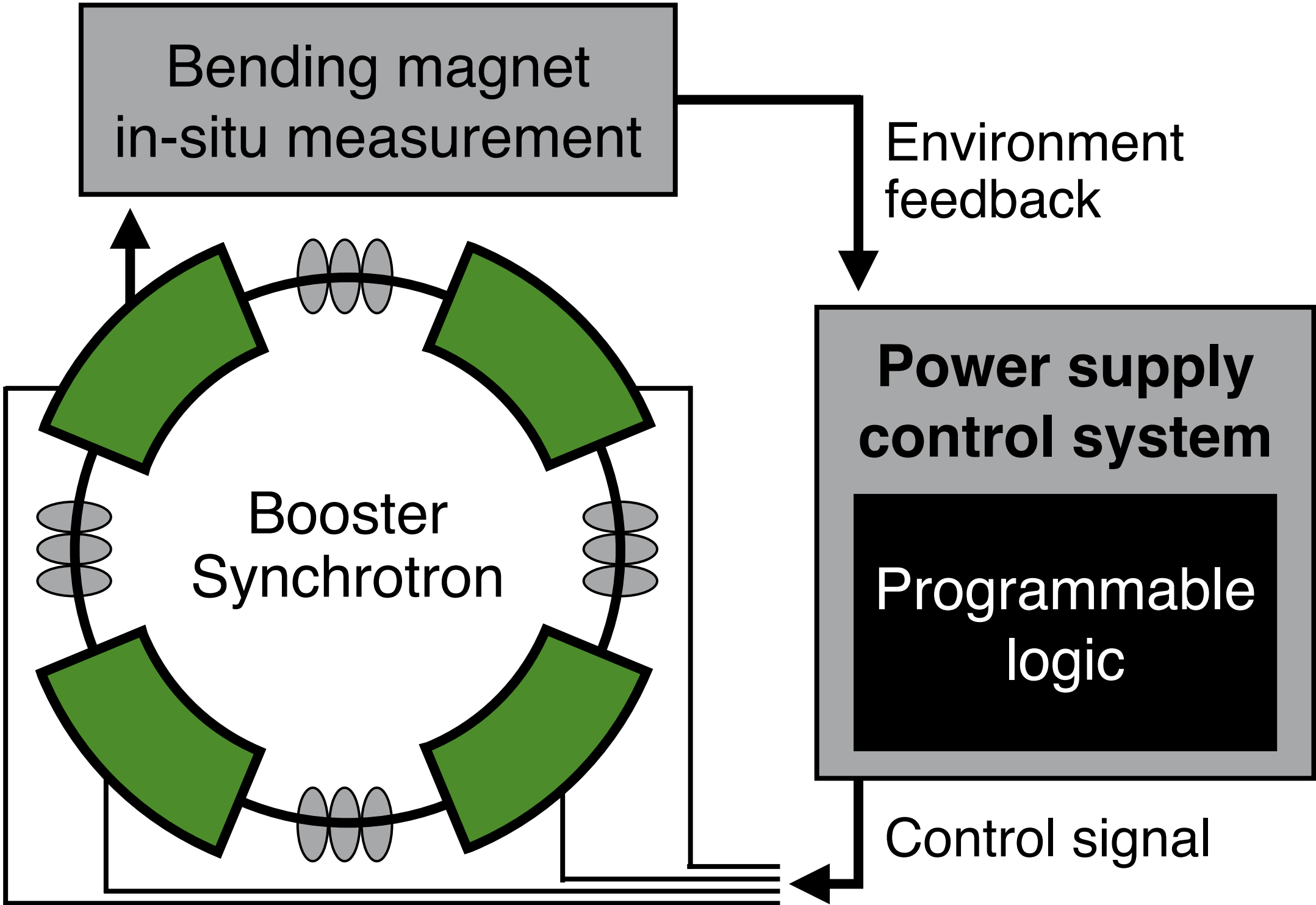| Type | Benchmark | Input Precision | Pipeline Rate | Real-time Latency | Misc. Req. | Baseline Model Parameters |
|---|---|---|---|---|---|---|
| Supervised Learning | Jet Classification | 16b | 150 ns | 1 $\mu$s | - | 4,389 |
| Unsupervised Learning | Sensor Data Compression | 9b | 25 ns | 100 ns | area, power (65 nm) | 2,288 |
| Reinforcement Learning | Beam Control | 32b | 5 ms | 5 ms | - | 34,695 |

| Type | Benchmark | Input Precision | Pipeline Rate | Real-time Latency | Misc. Req. | Baseline Model Parameters |
|---|---|---|---|---|---|---|
| Supervised Learning | Jet Classification | 16b | 150 ns | 1 $\mu$s | - | 4,389 |
| Unsupervised Learning | Sensor Data Compression | 9b | 25 ns | 100 ns | area, power (65 nm) | 2,288 |
| Reinforcement Learning | Beam Control | 32b | 5 ms | 5 ms | - | 34,695 |

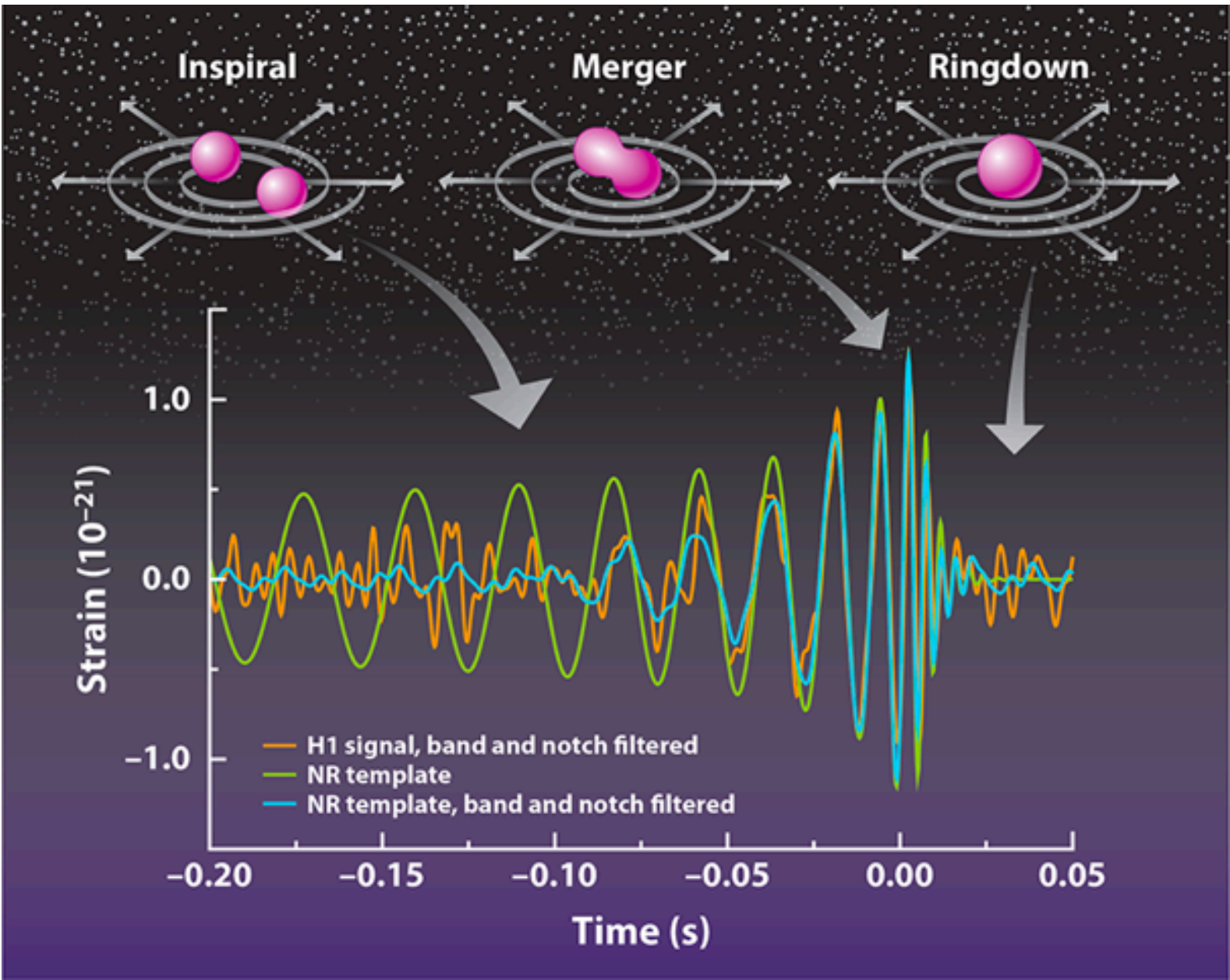▸ Particle jet classification for level-1 trigger: ~1 µs latency

| Type | Benchmark | Input Precision | Pipeline Rate | Real-time Latency | Misc. Req. | Baseline Model Parameters |
|---|---|---|---|---|---|---|
| Supervised Learning | Jet Classification | 16b | 150 ns | 1 $\mu$s | - | 4,389 |
| Unsupervised Learning | Sensor Data Compression | 9b | 25 ns | 100 ns | area, power (65 nm) | 2,288 |
| Reinforcement Learning | Beam Control | 32b | 5 ms | 5 ms | - | 34,695 |

▸ Particle jet classification for level-1 trigger: ~1 µs latency

▸ Sensor data compression: ~100 ns latency and additional area/power requirements
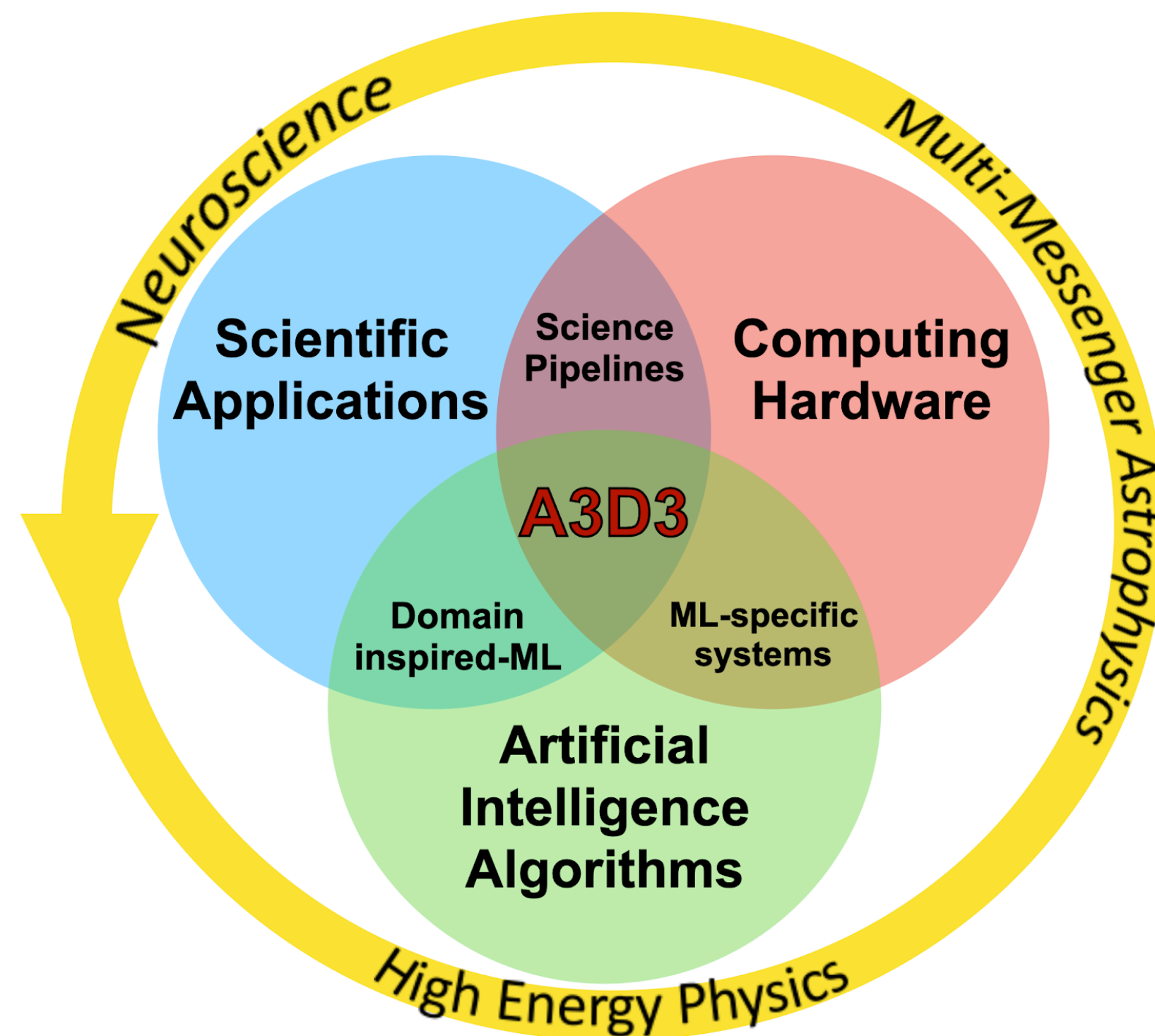
| Type | Benchmark | Input Precision | Pipeline Rate | Real-time Latency | Misc. Req. | Baseline Model Parameters |
|------|-----------|-----------------|---------------|-------------------|------------|---------------------------|
| Supervised Learning | Jet Classification | 16b | 150 ns | 1 $\mu$s | - | 4,389 |
| Unsupervised Learning | Sensor Data Compression | 9b | 25 ns | 100 ns | area, power (65 nm) | 2,288 |
| Reinforcement Learning | Beam Control | 32b | 5 ms | 5 ms | - | 34,695 |

‣ Particle jet classification for level-1 trigger: ~1 μs latency

‣ Sensor data compression: ~100 ns latency and additional area/power requirements

‣ Reinforcement learning for steering accelerator beams: ~5 ms latency

# CURRENT & FUTURE BENCHMARKS

| Type | Benchmark | Input Precision | Pipeline Rate | Real-time Latency | Misc. Req. | Baseline Model Parameters |
|---|---|---|---|---|---|---|
| Supervised Learning | Jet Classification | 16b | 150 ns | 1 $\mu$s | - | 4,389 |
| Unsupervised Learning | Sensor Data Compression | 9b | 25 ns | 100 ns | area, power (65 nm) | 2,288 |
| Reinforcement Learning | Beam Control | 32b | 5 ms | 5 ms | - | 34,695 |

▸ Particle jet classification for level-1 trigger: ~1 μs latency

▸ Sensor data compression: ~100 ns latency and additional area/power requirements

▸ Reinforcement learning for steering accelerator beams:  ~5 ms latency

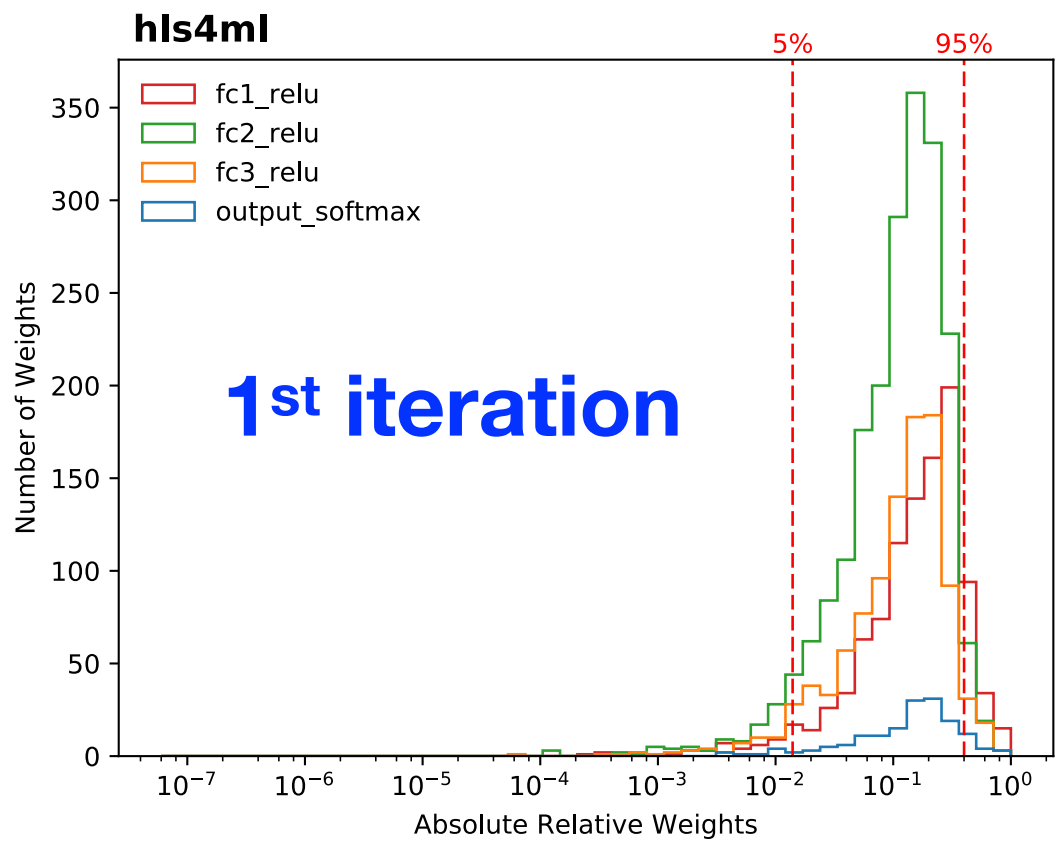▸ *Future: Time sequence analysis for gravitational wave or neural data, and more?*

▸ Tightly coupled organization of domain scientists, computer scientists, and engineers that unite three core components which are essential to achieve real-time AI to transform science: AI techniques, Computing Hardware, Scientific Applications
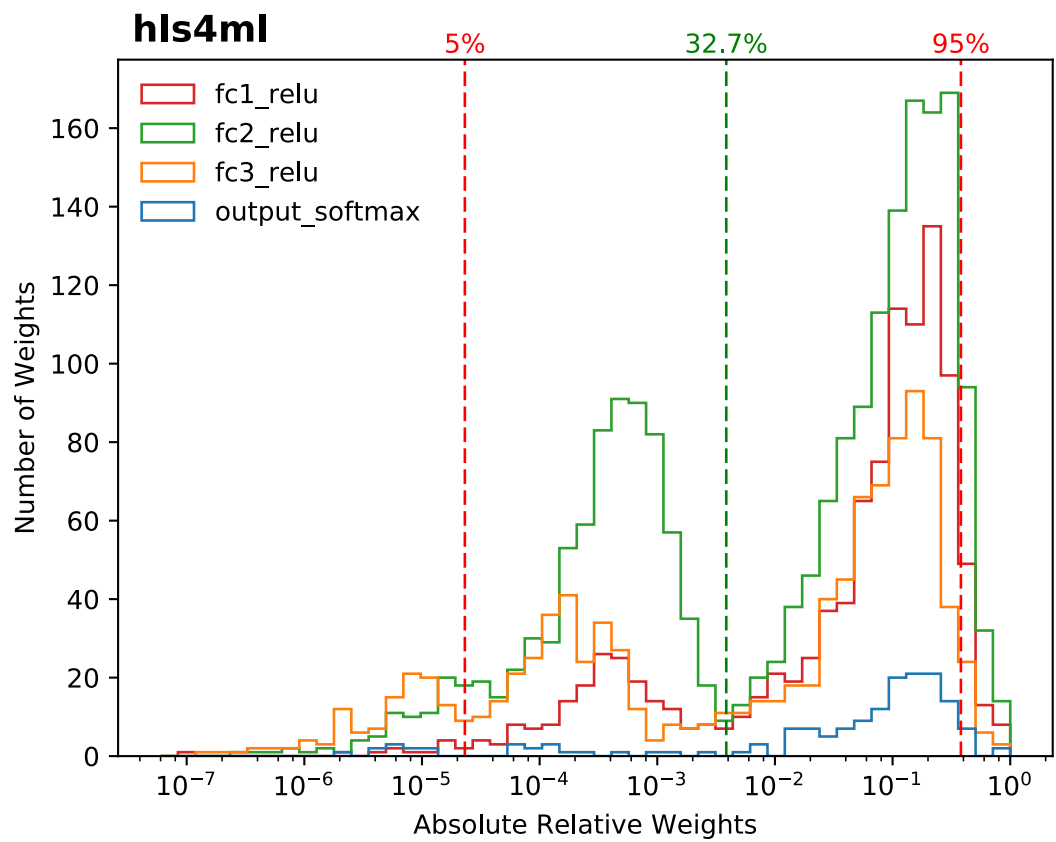
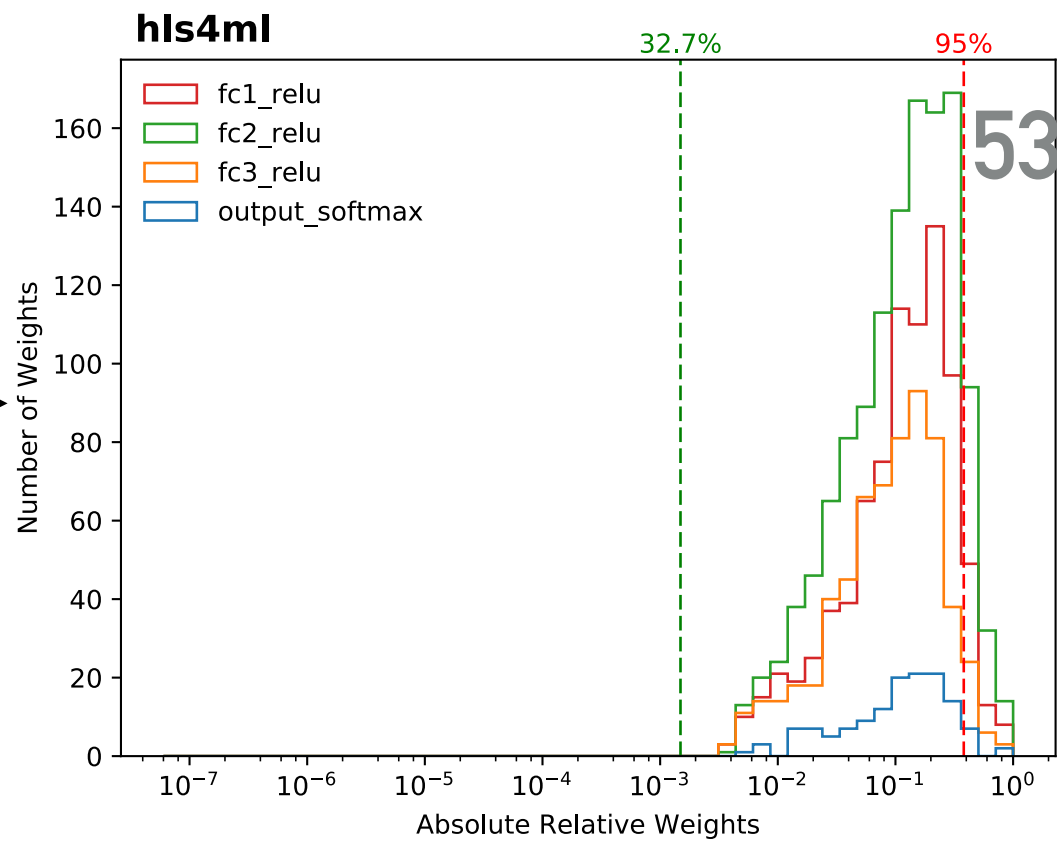▸ Collaborators welcome! Check the a3d3.ai for events

OAC-2117997